



# **European IT Certification Curriculum Self-Learning Preparatory Materials**

EITC/AI/MLP  
Machine Learning with Python



This document constitutes European IT Certification curriculum self-learning preparatory material for the EITC/AI/MLP Machine Learning with Python programme.

This self-learning preparatory material covers requirements of the corresponding EITC certification programme examination. It is intended to facilitate certification programme's participant learning and preparation towards the EITC/AI/MLP Machine Learning with Python programme examination. The knowledge contained within the material is sufficient to pass the corresponding EITC certification examination in regard to relevant curriculum parts. The document specifies the knowledge and skills that participants of the EITC/AI/MLP Machine Learning with Python certification programme should have in order to attain the corresponding EITC certificate.

#### Disclaimer

This document has been automatically generated and published based on the most recent updates of the EITC/AI/MLP Machine Learning with Python certification programme curriculum as published on its relevant webpage, accessible at:

<https://eitca.org/certification/eitc-ai-mlp-machine-learning-with-python/>

As such, despite every effort to make it complete and corresponding with the current EITC curriculum it may contain inaccuracies and incomplete sections, subject to ongoing updates and corrections directly on the EITC webpage. No warranty is given by EITCI as a publisher in regard to completeness of the information contained within the document and neither shall EITCI be responsible or liable for any errors, omissions, inaccuracies, losses or damages whatsoever arising by virtue of such information or any instructions or advice contained within this publication. Changes in the document may be made by EITCI at its own discretion and at any time without notice, to maintain relevance of the self-learning material with the most current EITC curriculum. The self-learning preparatory material is provided by EITCI free of charge and does not constitute the paid certification service, the costs of which cover examination, certification and verification procedures, as well as related infrastructures.

## TABLE OF CONTENTS

<b>Introduction</b>	<b>4</b>
Introduction to practical machine learning with Python	4
<b>Regression</b>	<b>6</b>
Introduction to regression	6
Regression features and labels	8
Regression training and testing	10
Regression forecasting and predicting	12
Pickling and scaling	14
Understanding regression	16
<b>Programming machine learning</b>	<b>17</b>
Programming the best fit slope	17
Programming the best fit line	19
R squared theory	20
Programming R squared	21
Testing assumptions	23
Introduction to classification with K nearest neighbors	26
K nearest neighbors application	28
Euclidean distance	31
Defining K nearest neighbors algorithm	33
Programming own K nearest neighbors algorithm	34
Applying own K nearest neighbors algorithm	36
Summary of K nearest neighbors algorithm	38
<b>Support vector machine</b>	<b>40</b>
Support vector machine introduction and application	40
Understanding vectors	41
Support vector assertion	43
Support vector machine fundamentals	45
Support vector machine optimization	47
Creating an SVM from scratch	53
SVM training	55
SVM optimization	57
Completing SVM from scratch	59
Kernels introduction	67
Reasons for kernels	69
Soft margin SVM	72
Soft margin SVM and kernels with CVXOPT	75
SVM parameters	77
<b>Clustering, k-means and mean shift</b>	<b>80</b>
Clustering introduction	80
Handling non-numerical data	84
K means with titanic dataset	86
Custom K means	88
K means from scratch	90
Mean shift introduction	92
Mean shift with titanic dataset	94
Mean shift from scratch	97
Mean shift dynamic bandwidth	100

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: INTRODUCTION****TOPIC: INTRODUCTION TO PRACTICAL MACHINE LEARNING WITH PYTHON**

Machine learning is a field of study that aims to give machines the ability to learn without being explicitly programmed to do so. In this programme, we will cover a variety of machine learning algorithms, including regression, classification with k-nearest neighbors and support vector machines, clustering with flat and hierarchical clustering, and deep learning with neural networks.

Each algorithm will be covered in three steps: theory, application, and inner workings. The theory provides high-level intuitions and can be quickly understood. Most algorithms are fairly basic to allow for scalability with large amounts of data. The application step involves using a module like scikit-learn to apply the algorithms to real-world data and observe their behavior. Finally, we will dive into the inner workings by recreating the algorithms from scratch in code, including all the math involved. This will provide a comprehensive understanding of how the algorithms work.

To follow along with this programme, it is recommended to have a basic understanding of Python 3. Additionally, a healthy amount of math will be covered, but it will be explained as we go along. Very basic familiarity with algebra and geometry is sufficient for most concepts.

Machine learning as a field has been around for over half a century, with Arthur Samuel defining it in 1959 as the study of giving machines the ability to learn without explicit programming. However, many people mistakenly believe that machine learning is hard-coded. In 1963, Vladimir Vapnik introduced the support vector machine, but it was not widely recognized until the 90s when it outperformed neural networks in handwritten character recognition. Support vector machines dominated the field until the recent resurgence of neural networks, particularly with deep learning.

If you feel like you are late to the machine learning party, rest assured that you are not. The computing power and accessibility we have today far surpass what was available in the past. In the 50s, computers could only handle a handful of bits at a time, and even in the 90s, it was challenging to access machines capable of running support vector machines at scale. However, now we have the ability to engage in deep learning with neural networks on gigabytes or even terabytes of data. With services like Amazon Web Services, we can easily rent GPU clusters for a fraction of the cost. We are living in an incredible time for machine learning.

Machine learning has become a popular field in artificial intelligence, but up until now, we have primarily focused on the learning aspect without involving the machine part. With tools like scikit-learn, you can achieve high accuracy without much understanding of the underlying algorithms. By simply applying default parameters, you can achieve around 90-95% accuracy. However, if you want to push the limits and obtain even higher accuracy, you need to delve deeper into how these algorithms work and learn how to tweak their parameters.

For instance, if you are working on a self-driving car, achieving 90-95% accuracy in distinguishing between a blob of tar and a child in a blanket is not sufficient. You need to aim for much higher accuracy. This material is designed for those who are eager to push the boundaries of what is currently available in machine learning.

The first topic we will cover is regression.

Regression is a technique used to predict continuous numerical values based on input variables. It is widely used in various fields, including finance, economics, and weather forecasting. In this material, we will explore different regression algorithms and learn how to implement them using Python.

To get started, we will discuss the fundamentals of regression, including the concepts of dependent and independent variables, as well as the different types of regression algorithms. We will then dive into linear regression, which is one of the most commonly used regression techniques. We will learn how to train a linear regression model, evaluate its performance, and make predictions.

Throughout this material, we will provide examples and code snippets to help you understand the concepts better. By the end, you will have a solid understanding of regression and be able to apply it to real-world

problems.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: REGRESSION****TOPIC: INTRODUCTION TO REGRESSION**

Regression is a technique used in machine learning to find the best fit line for continuous data. It involves modeling the relationship between features (attributes) and labels (continuous data) to predict future outcomes. In this example, we will use regression to analyze stock prices.

To get started, we need to install the necessary libraries. Open up the terminal or command prompt and install scikit-learn, Pandas, and Quandl by using the pip install command. Once installed, we can proceed with the example.

In regression, the equation of a straight line is used to model the data. The equation,  $y = mx + b$ , represents the dependent variable (y) as a function of the independent variable (x), with m as the slope and b as the y-intercept. The goal of regression is to determine the values of m and b that best fit the data.

In this case, we will be using the Quandl library to obtain stock price data. By specifying the ticker symbol, we can retrieve the dataset from Quandl. The dataset contains various features such as open, high, low, close, volume, and adjusted prices. These features represent different attributes of the stock.

When working with features, it is important to consider their relevance and meaningfulness to the data. Not all features may contribute significantly to the pattern recognition process. In this example, we will focus on the open, high, low, close, and adjusted prices, as they are closely related and provide meaningful data for our analysis.

Adjusted prices account for stock splits, where the number of shares and their prices are adjusted to maintain consistency. This ensures that the stock price does not appear to have changed drastically due to a split. We will be using the adjusted prices in our regression analysis.

It is worth noting that in regression, we do not consider the relationships between different features. However, in other machine learning algorithms such as deep learning, relationships between attributes can be explored. For regression, simplicity is key, and we aim to use meaningful features that have a direct impact on the data.

By simplifying our data and selecting relevant features, we can improve the accuracy and efficiency of our regression model. In the next steps, we will explore and analyze the dataset to further understand the relationship between the features and labels.

In this didactic material, we will introduce the concept of regression in the context of machine learning with Python. Regression is a supervised learning technique used to predict continuous numerical values based on input features. We will explore how to select relevant features and create a dataframe for regression analysis.

When working with machine learning classifiers, it is important to avoid including useless features as they can cause trouble. In supervised learning, especially with simpler classifiers, useless features can negatively impact the accuracy of predictions. Therefore, it is important to carefully choose the features that will be used in the analysis.

To begin, we will create a dataframe by selecting specific columns from an existing dataframe. The columns we will include are 'adjusted', 'open', 'high', 'low', 'close', and 'volume'. By selecting these columns, we are essentially recreating the dataframe to only include the relevant information needed for our analysis.

Some of these columns may seem relatively worthless at first glance, but they still hold valuable relationships. For example, the margin between the 'high' and 'low' prices can provide insights into the volatility of a given day. Similarly, the relationship between the 'open' and 'close' prices can indicate whether the price went up or down and by how much. These relationships can be valuable in predicting future outcomes.

However, a simple linear regression model will not automatically identify these relationships. It will only work with the features provided to it. Therefore, it is necessary to define these special relationships explicitly and use them as features in the regression analysis. By doing so, we can avoid redundant features that do not provide

additional useful information.

Let's start by calculating the 'high minus low percent', which represents the percent volatility. We will define a new column called 'HL\_percent', which is calculated as the difference between the 'high' and 'low' prices divided by the 'low' price, multiplied by 100. This calculation is performed on a per-row basis.

Next, we will calculate the 'percent\_change', which represents the daily percent movement. This is calculated similarly to the 'HL\_percent', but using the 'adjusted close' and 'adjusted open' prices instead. The 'percent\_change' is calculated as the difference between the 'adjusted close' and 'adjusted open' prices divided by the 'adjusted open' price, multiplied by 100.

Once we have calculated these new columns, we will define a new dataframe that only includes the columns we are interested in. In our case, the columns we care about are 'adjusted close', 'high low percent', 'percent change', and 'volume'. These columns provide us with the necessary information for our regression analysis.

Finally, we will print the first few rows of the dataframe to ensure that everything has been set up correctly. This step allows us to verify that the dataframe includes the desired columns and that the data is correctly formatted.

In the next tutorial, we will delve deeper into the concept of labels and features. We will explore whether the 'adjusted close' column will be used as a feature or as a label. This decision will depend on the specific goals of our analysis and the predictions we want to make.

If you have any questions or comments, please leave them below. Thank you for watching, and we appreciate your support and subscriptions. Stay tuned for the next tutorial, where we will continue our exploration of regression analysis with real data.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: REGRESSION****TOPIC: REGRESSION FEATURES AND LABELS**

In this tutorial, we will be discussing regression features and labels in the context of machine learning with Python. Regression is a type of supervised learning algorithm used for predicting continuous numerical values. Features are the input variables or attributes that are used to make predictions, while labels are the output variables that we want to predict.

In the previous tutorial, we discussed features and how to select them. Now, let's focus on defining the label. The label represents the target variable that we want to predict. In this case, we want to predict the future price of a stock. The only column we have that represents the price is the "adjusted close" column.

However, in order to predict the future price, we need to bring in some new information. We will be using the "adjusted close" column as a feature, but we also need to obtain the adjusted close in the future, maybe the next day or the next few days. This additional information will allow us to make accurate predictions.

To begin, we will define a variable called "forecast\_column" or "col", which will be equal to the "adjusted close" column. This variable can be changed to represent a different forecast column in the future, depending on the specific problem we are working on.

Next, we need to handle missing data. In machine learning, we cannot work with missing data, so we need to replace it with a specific value. In this case, we will replace missing data with -99,999. This value will be treated as an outlier in our dataset.

Now, let's discuss forecasting. Regression algorithms are commonly used for forecasting. In our case, we will define the variable "forecast\_out" as the integer value of  $\text{math.ceil}(0.1 * \text{length of the dataframe})$ . The  $\text{math.ceil}$  function rounds up any decimal number to the nearest integer. In this case, we are using 0.1 times the length of the dataframe to determine the number of days we want to forecast into the future.

It's important to note that the code provided in this tutorial is specific to stock prices, but the concepts can be applied to other regression problems as well. By changing the forecast column and adjusting the code accordingly, you can use similar techniques for different forecasting tasks.

In this tutorial, we discussed regression features and labels. We learned that features are the input variables used for making predictions, while labels are the output variables we want to predict. We also discussed the importance of handling missing data and how to forecast future values using regression algorithms.

In this didactic material, we will discuss the concept of regression features and labels in the context of machine learning with Python. Regression is a supervised learning algorithm used to predict continuous values based on input features. In this case, we will focus on predicting stock prices.

To begin, let's understand the process of creating regression features and labels. The features are the attributes that we believe may influence the target variable, which in this case is the adjusted close price of a stock. These features can be any relevant data points such as volume, previous prices, or other indicators.

In the provided transcript, the speaker mentions the need to forecast out a certain number of days. This means that we want to predict the stock price for a future time period. The number of days to forecast out is determined by the value of "forecast out" in the code. By default, the code sets this value to 0.1, which corresponds to 10% of the dataset.

To create the labels, we shift the adjusted close price column by the forecast out value. This means that each row's label will be the adjusted close price of the stock 10 days into the future. The code achieves this by using the "shift" function in pandas.

It is important to note that the forecast out value can be changed to suit different prediction timeframes. For example, if we want to predict the stock price for tomorrow only, we can set the forecast out value to 0.01.



Once the features and labels are created, we can proceed with training and testing our regression model. However, it is worth mentioning that regression alone may not lead to substantial financial gains. It serves as a good starting point for modeling stock prices, but additional useful features should be incorporated to improve the accuracy of predictions.

At this point, we have covered the process of creating regression features and labels for predicting stock prices using Python. In the next material, we will consider training, testing, and running the regression algorithm on real data.

If you have any questions or concerns, please feel free to leave them in the comments section. Stay tuned for the next material where we will explore the practical application of regression in predicting stock prices.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: REGRESSION****TOPIC: REGRESSION TRAINING AND TESTING**

In this tutorial, we will focus on regression training and testing in the context of machine learning with Python. Regression is a supervised learning technique used to predict continuous numerical values based on input features. We will be using the scikit-learn library, which provides a wide range of machine learning algorithms and tools.

Before we begin, let's make some necessary imports. We will import numpy as NP and scikit-learn's preprocessing module. Numpy is a powerful computing library that allows us to work with arrays, while preprocessing provides various data scaling and transformation methods. Additionally, we will import cross-validation and support vector machine (SVM) modules, which we will use later in the tutorial.

To start, we need to define our features and labels. Features, denoted as  $x$ , represent the input variables, while labels, denoted as  $y$ , represent the output variable we want to predict. We will use a numpy array to store the features, obtained by dropping the label column from our dataset. Similarly, we will store the labels in a numpy array.

Next, we will scale our features using the `preprocessing.scale` function. Scaling the data is typically done to ensure that the features are within a specific range, often between -1 and 1. This can improve accuracy and processing speed. It's important to note that if we plan to use the classifier in real-time on new data, we need to scale the new values alongside the training data.

After scaling, we will redefine our features array,  $x$ , to include only the points for which we have corresponding labels,  $y$ . This ensures that our dataset is aligned properly. We will also drop any rows with missing values using the `df.dropna` method.

Finally, we will print the lengths of  $x$  and  $y$  to verify that we have the correct number of data points.

We have covered the initial steps of regression training and testing. We imported the necessary libraries, defined our features and labels, scaled the features, and ensured the alignment of our dataset. In the next tutorial, we will continue with the regression process.

In the process of regression training and testing in machine learning with Python, it is important to create training and testing sets. To do this, we can use the `train_test_split` function from the scikit-learn library. The function takes in the features ( $x$ ) and labels ( $y$ ), as well as the desired test size (in this case, 20% or 0.2). The function shuffles the data while keeping the features and labels connected, and outputs the training and testing data for both  $x$  and  $y$ .

Once we have the training data, we can proceed to fitting a classifier. In this example, we will use linear regression as the classifier. To fit the classifier, we use the `fit()` function and pass in the features ( $x_{\text{train}}$ ) and labels ( $y_{\text{train}}$ ). This trains the classifier on the training data.

After fitting the classifier, it is important to test its performance. This is done by using the `score()` function, which is synonymous with testing. We pass in the testing features ( $x_{\text{test}}$ ) and labels ( $y_{\text{test}}$ ) to get the accuracy score. It is worth noting that training and testing on separate data is important to avoid overfitting, where the classifier simply memorizes the training data and performs poorly on new data.

In this example, the accuracy score obtained is 0.96, indicating a 96% accuracy in predicting the price with a one-day shift. The accuracy score is calculated using squared error, which will be explained in detail in the upcoming explanation of linear regression.

It is also possible to use other algorithms besides linear regression. For example, support vector regression (SVR) can be used. Switching to SVR is as simple as changing the classifier to SVR in the code. However, it is important to note that the performance of different algorithms can vary significantly. In this example, SVR performs much worse than linear regression, with an accuracy score of only 0.51. It is worth experimenting with different algorithms and their parameters, such as different kernels, to improve performance.

Regression training and testing in machine learning with Python involves creating training and testing sets, fitting a classifier, and evaluating its performance using accuracy scores. It is important to choose the right algorithm and parameters to achieve accurate predictions.

Support Vector Regression is not the focus of this tutorial series. However, it is important to understand the concept of a kernel, which will be explained in the support vector machines tutorial. This example demonstrates how easy it is to switch between different algorithms for regression, classification, or clustering. It is recommended to test different algorithms to find the most suitable one for your task.

When working with various algorithms, it is essential to refer to the documentation. For instance, when using linear regression, you should check if the algorithm supports threading. Threading allows running multiple jobs or threads simultaneously. Linear regression can be threaded easily, unlike support vector machines, which require more complex techniques. By threading linear regression, you can perform parallel operations, resulting in faster training.

To determine if an algorithm can be threaded, you can search for it on Google and look for the "N Jobs" parameter. This parameter indicates the number of jobs or threads that can be run concurrently. By default, linear regression runs with only one job. However, you can specify a higher number to run multiple jobs simultaneously and increase the speed of training. Alternatively, you can use "-1" to utilize as many jobs as possible based on your processor's capabilities.

It is important to consider your computer's processing power while following this tutorial series. If you have an older computer or a laptop, some operations may take longer to run. In such cases, it might be beneficial to use a server or more powerful hardware, especially when working with deep learning algorithms.

Understanding the threading capabilities of different algorithms is important. Linear regression, as well as many other algorithms, can be highly threaded, allowing for efficient parallel operations. As processing power increases, the ability to scale calculations and methodologies becomes increasingly valuable.

In the next tutorial, we will discuss predicting future values using Scikit-learn. Following that, we will consider linear regression, providing a detailed breakdown and implementation. If you have any questions or concerns, feel free to leave them in the comments section. Thank you for your support and stay tuned for the upcoming tutorials.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: REGRESSION****TOPIC: REGRESSION FORECASTING AND PREDICTING**

In this tutorial, we will be building on the previous regression tutorial and focus on regression forecasting and predicting using machine learning with Python.

To begin, we have already created a linear regression algorithm in the previous tutorial and found that it has great accuracy. Now, we are ready to predict values beyond our known data.

We already have some unknown data because we are forecasting into the future, specifically for a period of 30 days. To work with this unknown data, we need to modify our X's. We will define a new variable called `x_lately`, which will contain the X values for the unknown period.

To do this, we will set `x_lately` equal to `x[-forecast_out:]`, where `forecast_out` is the number of days we want to forecast. This will give us the X values for the unknown period.

Next, we need to figure out the values of M and B in the equation  $y = Mx + B$ . We have already done linear regression to find these values for the known data. Now, we will use these values to predict the values for the unknown data.

To make predictions, we need to create a forecast set. We will use the classifier's `predict` method to make predictions based on the `x_lately` data. We can pass a single value or an array of values to predict multiple values at once. In this case, we will pass the `x_lately` array to predict the values for the entire unknown period.

Once we have the forecast set, we can print the predicted values, the confidence (or accuracy) of the predictions, and the number of days we are forecasting out.

Finally, if we want to graph the predicted values, we can use the `matplotlib` library. We will import `datetime` and `matplotlib.pyplot` to help us with graphing. We will also import the `style` module from `matplotlib` to specify the style of the graph.

To create the graph, we will create a new column in our data frame called `df_forecast` and fill it with NaN values. This column will be used to plot the predicted values. We will also find the last date in our data to use as a reference for the graph.

With all the necessary preparations, we can now plot the graph using the specified style. This will give us a visual representation of the predicted values for the next 30 days.

Regression forecasting and predicting is an important aspect of machine learning. In this process, we use historical data to make predictions about future outcomes. In this didactic material, we will discuss how to perform regression forecasting and predicting using Python.

To begin with, let's understand the concept of regression. Regression is a statistical technique used to model the relationship between a dependent variable (also known as the label) and one or more independent variables (also known as features). In the context of machine learning, regression is used for predicting continuous numerical values.

In the given material, the speaker discusses the process of creating a graph to visualize the forecasted data. The speaker mentions that the dates are not included as features in the regression model, but they are needed for plotting the graph. To address this, the speaker uses the `'df.loc'` function to reference the index of the DataFrame, which represents the dates. By doing this, the speaker ensures that the dates are included on the axes of the graph.

The speaker then proceeds to populate the DataFrame with the new dates and forecasted values. This is done using a for loop, where each forecasted value and date are iteratively added to the DataFrame. The speaker also mentions that the `'df.loc'` function creates the index if it doesn't exist and replaces it if it does.

Once the DataFrame is populated, the speaker plots the actual data and the forecasted data on a graph using the 'df.Adj.Close.plot' and 'df.Forecast.plot' functions, respectively. The 'plt.legend' function is used to add a legend to the graph, and the 'plt.xlabel' and 'plt.ylabel' functions are used to label the axes.

The speaker concludes by mentioning that the purpose of including the dates in the graph is to provide a visual representation of the forecasted data. The speaker also highlights that the complex part of the process is to ensure that the dates are included on the axes, which requires some additional steps.

This didactic material explains the process of regression forecasting and predicting using Python. It covers the steps involved in creating a graph to visualize the forecasted data and highlights the importance of including the dates on the axes of the graph.

In this didactic material, we will discuss regression forecasting and predicting in the context of machine learning with Python. Regression is a supervised learning algorithm used to predict continuous output variables based on input features. It is commonly used for tasks such as forecasting stock prices, predicting house prices, or estimating sales figures.

To begin, let's assume we have a dataset with multiple columns of input features and one column of output values that we want to predict. We start by creating a regression model using Python libraries such as scikit-learn. The model learns the relationship between the input features and the output variable from the training data.

Once we have trained the regression model, we can use it to make predictions on new data. In the context of forecasting, we often want to predict future values based on historical data. To achieve this, we can add forecasts at the end of our dataset. This is a simple and practical approach, although it may not be the most elegant solution.

One important concept to mention is the use of 'pickling' in machine learning. Pickling refers to the process of saving a trained model to a file, which allows us to reuse the model without the need for retraining. This can be particularly useful when dealing with large datasets or when we need to make frequent predictions.

For example, if we have a classifier trained on a relatively small dataset, it may take a long time to train the model every time we want to make a prediction. By pickling the model, we can quickly load it in without any training time, making the prediction process much more efficient.

Regression forecasting and predicting are important techniques in machine learning. Regression models allow us to predict continuous output variables based on input features. By pickling our trained models, we can save time and computational resources when making predictions on new data.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: REGRESSION****TOPIC: PICKLING AND SCALING**

In this tutorial, we will discuss the concepts of pickling and scaling in the context of machine learning with Python. Pickling refers to the serialization of any Python object, such as a dictionary or a classifier. It allows us to save objects so that we can reuse them later, saving time and effort. Scaling, on the other hand, involves transforming the input features of a dataset to a specific range. This is done to ensure that all features have a similar scale and to prevent any one feature from dominating the learning process.

To begin, we need to import the 'pickle' module. This module provides functions for pickling and unpickling objects. We can import it using the statement 'import pickle'. Next, we will explore how pickling works. Think of pickling as saving a file. We open the file, save it, and then when we want to use it, we open and read it. In the case of pickling, we open a file with the intention to write, and then we use the 'pickle.dump()' function to save the trained classifier object. The classifier can be any machine learning model that we have trained. We specify the file name and the mode as 'wb' (write binary) to ensure that the file is saved in binary format.

To use the classifier, we need to unpickle it. We open the file in read binary mode using the 'pickle.load()' function and assign the loaded classifier to a variable. This allows us to make predictions without having to retrain the model each time. It is important to note that pickling is useful when we have large amounts of data and training the model is time-consuming. By saving the trained model, we can avoid the training step and directly use the classifier for predictions.

It is worth mentioning that in order to use pickling effectively, it is recommended to retrain the classifier periodically, such as once a month. This ensures that the model remains up-to-date with any changes in the data.

Scaling, on the other hand, involves transforming the input features of a dataset to a specific range. This is important because features with different scales can have a significant impact on the performance of machine learning algorithms. Scaling ensures that all features have a similar scale, preventing any one feature from dominating the learning process. There are various scaling techniques available, such as Min-Max scaling and Standardization. These techniques can be applied using libraries like 'scikit-learn'.

Pickling allows us to save trained machine learning models, such as classifiers, so that we can reuse them later without having to retrain the models. This saves time and effort, especially when dealing with large datasets. Scaling, on the other hand, ensures that all input features have a similar scale, preventing any one feature from dominating the learning process. By scaling the features, we can improve the performance of machine learning algorithms.

In the previous material, we discussed the concepts of pickling and scaling in the context of machine learning with Python. Pickling refers to the process of saving a trained model to a file, while scaling involves transforming the input features to a specific range. These techniques are commonly used in machine learning workflows to improve efficiency and accuracy.

When we pickle a model, we serialize it into a file that can be easily stored and retrieved later. This allows us to reuse the trained model without having to retrain it every time. In Python, the pickle module provides functions for pickling and unpickling objects. To pickle a classifier, we simply need to call the pickle.dump() function and pass the classifier object and the file object as arguments. Conversely, to unpickle a classifier, we use the pickle.load() function and pass the file object as an argument.

Scaling, on the other hand, involves transforming the input features to a specific range. This is important because features with different scales can have a disproportionate impact on the model's performance. One common scaling technique is standardization, which transforms the features to have zero mean and unit variance. In Python, the scikit-learn library provides a StandardScaler class for scaling the features. We can create an instance of the StandardScaler class and use its fit\_transform() method to scale the input features.

It is worth noting that linear regression algorithms can be scaled effectively. This means that scaling the input features can have a positive impact on the performance of linear regression models. Therefore, it is

recommended to scale the features before training a linear regression model.

In the next material, we will dive deeper into the topic of linear regression and learn how to implement our own linear regression algorithm. This will provide us with a better understanding of how linear regression works and how it can be applied to real-world problems.

If you have any questions or comments, please feel free to leave them below. Thank you for watching and for your continued support. Stay tuned for more exciting material on machine learning with Python!

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: REGRESSION****TOPIC: UNDERSTANDING REGRESSION**

Linear regression is a fundamental concept in machine learning that involves predicting a continuous output variable based on one or more input variables. In this didactic material, we will discuss linear regression and its implementation in Python.

Linear regression is used to model the relationship between a dependent variable (y) and one or more independent variables (x). The goal is to find the best-fit line that minimizes the difference between the predicted values and the actual values of the dependent variable. This line is represented by the equation  $y = mx + b$ , where m is the slope and b is the y-intercept.

To understand linear regression, let's consider some examples. In the first example, we have a dataset with data points that form a straight line. This indicates a strong positive correlation between the independent and dependent variables. In the second example, the data points form a curve, indicating a weaker correlation. Finally, in the third example, the data points do not show any clear relationship. In this case, linear regression would not be beneficial.

To calculate the slope (m) and y-intercept (b) for the best-fit line, we use the following formulas:

$$m = (\text{mean}(x) * \text{mean}(y) - \text{mean}(x * y)) / (\text{mean}(x)^2 - \text{mean}(x^2))$$

$$b = \text{mean}(y) - m * \text{mean}(x)$$

Here, mean(x) represents the mean of all x values, mean(y) represents the mean of all y values, and x\*y represents the product of each x and y value. The formulas for m and b allow us to define the equation of the best-fit line.

Once we have calculated the values of m and b, we can use them in the equation  $y = mx + b$  to predict the values of y for any given x. By plugging in different x values, we can generate a line that represents the relationship between the independent and dependent variables.

In Python, we can implement linear regression using pure code. By using libraries such as NumPy and Matplotlib, we can easily perform calculations and visualize the results. By understanding the theory behind linear regression and implementing it in Python, we can gain insights into how this algorithm works and how it can be applied to real-world datasets.

Linear regression is a powerful technique for predicting continuous output variables based on input variables. By understanding the concepts of slope and y-intercept, we can calculate the best-fit line and make predictions. Python provides tools and libraries that make it easy to implement linear regression and visualize the results.

Regression is an important concept in machine learning that involves finding the best-fit line given a set of x's and y's. Although the math behind it can become more complex as the dimensions in vector space increase, this example focuses on a simple regression problem. The algorithm used to find the values for the best-fit line is the key component in regression.

To apply regression in Python, we can translate the algorithm into code. In the next tutorial, we will convert the algorithm into Python code and apply it to real data. This will allow us to see regression in action and gain a better understanding of its practical applications.

If you have any questions, comments, or concerns, please leave them below. Thank you for watching and for your continued support and subscriptions. Stay tuned for the next tutorial!



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: PROGRAMMING THE BEST FIT SLOPE**

In this part of our machine learning tutorial series, we will be focusing on creating a simple linear regression algorithm from scratch using Python. The goal is to understand how to program the best fit slope for a given dataset.

To begin, let's recall the equation for a line:  $y = MX + B$ . In this equation,  $X$  represents the x-axis values, while  $M$  and  $B$  are the slope and y-intercept, respectively. Our first step is to calculate the slope,  $M$ .

The formula for calculating the slope is as follows:

$$M = ((\text{mean of } X \text{ values} * \text{mean of } Y \text{ values}) - (\text{mean of } X \text{ values} * Y \text{ values})) / ((\text{mean of } X \text{ values})^2 - \text{mean of } (X \text{ values}^2))$$

To translate this formula into Python code, we need to import the mean function from the statistics module, as well as the numpy module as NP. We will use the mean function to calculate the means of the  $X$  and  $Y$  values.

Next, we define some simple values for the  $X$  and  $Y$  variables. For example,  $X = [1, 2, 3, 4, 5, 6]$  and  $Y = [5, 4, 6, 5, 6, 7]$ . We can visualize this data using the matplotlib module, which we import as PLT.

After visualizing the data, we convert the  $X$  and  $Y$  variables to numpy arrays using the `numpy.array()` function. We also specify the data type as float64 using `NP.float64`.

Now, let's define a function to calculate the best fit slope. We will pass the  $X$  and  $Y$  variables as arguments and return the slope,  $M$ .

The first step in calculating the slope is to find the mean of the  $X$  values multiplied by the mean of the  $Y$  values. We store this value in the variable  $M$ .

Next, we subtract the mean of the  $X$  values multiplied by the  $Y$  values from  $M$ .

Finally, we divide this result by the difference between the mean of the  $X$  values squared and the mean of the  $X$  values squared.

At this point, we have completed the top part of the fraction in the formula.

To continue, we add another set of parentheses and subtract the mean of the  $X$  values multiplied by the  $Y$  values.

Returning to our function, we have now completed the entire top part of the fraction.

The next step is to add a third set of parentheses in the code.

Please note that the code provided here is a simplified version for educational purposes. In practice, you may need to handle additional complexities and data preprocessing steps before calculating the best fit slope.

In this tutorial, we will continue our discussion on programming machine learning with Python, specifically focusing on programming the best fit slope.

To begin, let's recap what we have learned so far. In the previous tutorial, we discussed the concept of mean and how to calculate it using Python. Now, we will explore how to calculate the power of two, denoted by  $^2$ , which is essentially the mean of the  $X$ 's multiplied by the mean of the  $X$ 's.

In Python, there are a few different ways to calculate the power of two. One option is to use the caret symbol ( $^$ ), which represents exponentiation. However, when we run this code, we may encounter an error message stating "unsupported operand for the data type we're using."

---

**EUROPEAN IT CERTIFICATION CURRICULUM SELF-LEARNING PREPARATORY MATERIALS**

---

Another option is to use the asterisk symbol (\*) for multiplication. This method is acceptable and will give us the desired result. Alternatively, we can also use the explicit multiplication notation, such as "mean of the X's times the mean of the X's." Both of these approaches will yield the same outcome.

Moving on, we need to calculate the mean of the X's squared, denoted by  $X^2$ . To accomplish this, we can use the same multiplication notation as before. However, we need to subtract the mean of the X's from this value. Again, we have a few options to achieve this. We can use the caret symbol (^) or the asterisk symbol (\*) for multiplication.

Once we have obtained the desired values, we can proceed to calculate the best fit slope, denoted by M. We can subtract the mean of the X's squared from the mean of the X's, and then divide the result by the mean of the X's squared minus the X's squared. It is important to note that the order of operations, known as PEMDAS (parentheses, exponents, multiplication, division, addition, subtraction), must be followed to obtain the correct result.

In some cases, we may encounter issues with the order of operations if we do not use parentheses correctly. For example, if we divide before subtracting, we may obtain unexpected results. Therefore, it is important to use parentheses appropriately to ensure the desired outcome.

After calculating the best fit slope, we obtain a value of -15.26. It is worth mentioning that a negative slope is unusual in the context of positively correlated data. However, this is just an example to illustrate the programming process.

Finally, we need to calculate the intercept, denoted by B. This will be the focus of our next tutorial, where we will discuss linear regression. Stay tuned for the next video! If you have any questions or comments, please feel free to leave them below. Thank you for watching and for your continued support.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: PROGRAMMING THE BEST FIT LINE**

In this tutorial, we will continue our exploration of linear regression in machine learning. Specifically, we will focus on calculating the y-intercept of the best fit line.

To recap, the best fit line is calculated using the slope (M) and the y-intercept (B). The equation for the y-intercept of the best fit line is  $B = \text{mean}(Y) - M * \text{mean}(X)$ . Once we have the slope, calculating the y-intercept is straightforward.

In our function, we will modify it to include both the best fit slope and intercept. We will define M as before and also return B. To calculate B, we simply set B equal to the mean of the Y values minus M times the mean of the X values. Finally, we will return both M and B.

To create a line that fits the data, we can use the equation  $y = MX + B$ . We can generate a list of Y values using a one-line for loop. For each X value in our dataset, we calculate  $MX + B$ . This creates a regression line that represents the best fit for our data.

To visualize the regression line, we will use the matplotlib library. We import the style module and set it to '538'. Next, we scatter plot the X and Y values and plot the regression line using the X values and the calculated regression line. Finally, we add a title to the plot.

Now, let's discuss the significance of the regression line. By creating a model for the data using the equation  $y = MX + B$ , we can make predictions based on the model. For example, if we want to predict the Y value when X equals 8, we can simply calculate  $MY + B$ . This allows us to make predictions based on our model. We can even visualize the prediction by scatter plotting it with a green color.

At this point, we have successfully created a model for our data and made predictions based on that model. However, it is important to assess the accuracy of our best fit line. We want to ensure that it is not only the best fit line but also a good fit line for our data. To evaluate the accuracy, we can calculate how well the best fit line fits the data. This will give us an indication of the line's accuracy.

In this tutorial, we learned how to calculate the y-intercept of the best fit line and how to create a regression line that fits the data. We also explored how to make predictions based on the model and assessed the accuracy of the best fit line. By understanding these concepts, we can effectively use linear regression in machine learning applications.

In the field of machine learning, one important concept is linear regression. Linear regression involves finding the best fit line for a given set of data points. This line represents the relationship between the independent variable(s) and the dependent variable. The accuracy and confidence of the best fit line are important factors in evaluating its performance.

Accuracy and confidence are closely related in linear regression, especially when dealing with a small number of dimensions. The accuracy of the best fit line refers to how well it predicts the dependent variable based on the independent variable(s). It measures the closeness of the predicted values to the actual values. On the other hand, confidence in this context refers to the certainty or reliability of the predictions made by the best fit line. It indicates how much trust we can place in the predictions.

However, as the number of dimensions increases, such as in the case of K nearest neighbors, the difference between accuracy and confidence becomes more pronounced. In K nearest neighbors, accuracy refers to how well the algorithm classifies new data points based on their proximity to the training data. Confidence, on the other hand, reflects the algorithm's certainty in its predictions.

Moving forward, the next step in our journey is to determine how good of a fit the best fit line is. This will be the focus of our upcoming tutorial. If you have any questions or comments, please feel free to leave them below. Otherwise, stay tuned for the next material.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: R SQUARED THEORY**

Welcome to this section of our machine learning tutorial series, where we will discuss the concept of R-squared or the coefficient of determination in the context of linear regression. After calculating the best fit line in our Python code, it is important to determine the accuracy of this line. This is where R-squared comes into play.

R-squared is calculated using squared error, which measures the accuracy of the best fit line. To understand squared error, let's consider an example. Imagine you have two graphs with some data points, and you want to draw the best fit line. One graph might have data points closer to the line, while the other might have points further away. In this case, we would consider the graph with points closer to the line as a better fit.

To calculate squared error, we measure the distance between each data point and the best fit line, and then square that value. By squaring the error, we ensure that we are only dealing with positive values. Additionally, squaring the error helps penalize for outliers, as we want to focus on linear regression for linear data.

Now, let's move on to calculating R-squared. R-squared is obtained by subtracting the squared error of the best fit line from 1, and then dividing it by the squared error of the mean of the Y values in the dataset. In other words, we compare the accuracy of the best fit line to the accuracy of a simple straight line representing the mean of the Y values.

A good value for R-squared indicates a strong fit between the best fit line and the data, while a bad value suggests a weak fit. For example, an R-squared value of 0.8 means that the squared error of the best fit line is 20% of the squared error of the mean of the Y values.

R-squared is a useful metric for evaluating the accuracy of a best fit line in linear regression. By calculating the squared error and comparing it to the squared error of the mean of the Y values, we can determine how well our model fits the data.

In machine learning, the concept of R squared theory is used to evaluate the performance of a model. R squared measures the proportion of the variance in the dependent variable that can be explained by the independent variable(s). A high R squared value indicates a good fit of the model to the data.

To calculate R squared, we first need to calculate the squared error. The squared error is the square of the difference between the predicted values ( $\hat{Y}$ ) and the actual values ( $Y$ ). We compare this squared error to the squared error of the mean of the Y values.

For example, let's say the squared error of the  $\hat{Y}$  line is 2 and the squared error of the mean of the Y values is 10. In this case, the squared error of the  $\hat{Y}$  line is significantly lower than the squared error of the mean of the Y values. This indicates that the model is performing well and the data is likely linear. An R squared value of 0.8 is considered good.

On the other hand, if the R squared value is low, such as 0.3, it means that the squared error of the  $\hat{Y}$  line is closer to the squared error of the mean of the Y values. This indicates a poorer fit of the model to the data.

To calculate R squared, we divide the squared error of the  $\hat{Y}$  line by the square root of the mean of the Y values. We want this value to be as close to 0 as possible. The higher the R squared value, the better the accuracy of the model.

It's important to note that R squared is not a percent accuracy, but rather a coefficient of determination. It measures the proportion of the variance explained by the model.

In Python, we can easily calculate R squared using the squared error and mean of the Y values. This allows us to evaluate the performance of our machine learning models.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: PROGRAMMING R SQUARED**

Hello and welcome to this machine learning tutorial. In this material, we will be building upon our previous knowledge on calculating the coefficient of determination, also known as the r-squared value. The r-squared value is an important measure of how well our best fit line fits the data.

To calculate the r-squared value, we first need to understand the concept of squared error. Squared error is the difference between the original y-values and the y-values predicted by the line. It represents the amount of error in the y-values, which is then squared. We can define a function called "squared error" to calculate this.

To calculate the squared error for the entire line, we sum up the squared differences between the predicted y-values and the original y-values. This can be done by subtracting the y-values predicted by the line from the original y-values, squaring the result, and summing up all these squared differences.

Once we have the squared error, we can proceed to calculate the coefficient of determination. The coefficient of determination is calculated as 1 minus the squared error of the line divided by the squared error of the mean of the y-values. It represents the proportion of the variance in the y-values that is explained by the line.

To calculate the coefficient of determination, we first need to calculate the mean of the y-values. This can be done by taking the sum of all the original y-values and dividing it by the total number of y-values.

Next, we calculate the squared error of the line by using the squared error function we defined earlier. We pass in the original y-values and the y-values predicted by the line.

Finally, we calculate the coefficient of determination by subtracting the squared error of the line from 1 and dividing it by the squared error of the mean line.

In our case, the coefficient of determination is 0.58. This means that the line we are analyzing explains 58% of the variance in the y-values. A coefficient of determination of 0 would indicate that the line is as accurate as the mean of the y-values, while a coefficient of determination of 1 would indicate a perfect fit.

It is important to note that the squared error and coefficient of determination are not the only measures of accuracy for a best fit line. However, they provide valuable insights into how well the line fits the data.

In the next tutorial, we will be discussing the testing of our assumptions and using sample data to validate our algorithms. This will help us ensure the accuracy of our calculations and identify any potential errors.

Artificial Intelligence (AI) and Machine Learning (ML) have become integral parts of many industries, enabling computers to learn and make predictions or decisions without explicit programming. In this didactic material, we will focus on programming machine learning algorithms using Python, specifically exploring the concept of R squared.

Python is a popular programming language for ML due to its simplicity and extensive libraries. To begin programming ML algorithms, we need to ensure that Python and the necessary libraries are installed on our system. Once installed, we can import the required libraries, such as NumPy and scikit-learn, which provide powerful tools for scientific computing and ML.

A fundamental concept in ML is supervised learning, where a model learns from labeled training data to make predictions on unseen data. One commonly used metric to evaluate the performance of regression models is R squared ( $R^2$ ).  $R^2$  measures the proportion of the variance in the dependent variable that can be explained by the independent variables.

To calculate  $R^2$ , we need to understand the concept of variance. Variance measures the spread or dispersion of a set of values. In ML, we often encounter the terms "explained variance" and "total variance." Explained variance refers to the variance explained by the model, while total variance represents the overall variance in the dependent variable.

The formula to calculate  $R^2$  is as follows:

$$R^2 = 1 - (\text{explained variance} / \text{total variance})$$

In Python, we can calculate  $R^2$  using the scikit-learn library. First, we need to split our dataset into training and testing sets. We then fit our model on the training data and make predictions on the testing data. Finally, we can calculate  $R^2$  using the `score` method provided by scikit-learn.

Here is an example code snippet showcasing the calculation of  $R^2$  using scikit-learn:

1.	from sklearn.model_selection import train_test_split
2.	from sklearn.linear_model import LinearRegression
3.	
4.	# Split the dataset into training and testing sets
5.	X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
6.	
7.	# Create a Linear Regression model
8.	model = LinearRegression()
9.	
10.	# Fit the model on the training data
11.	model.fit(X_train, y_train)
12.	
13.	# Make predictions on the testing data
14.	y_pred = model.predict(X_test)
15.	
16.	# Calculate $R^2$
17.	r_squared = model.score(X_test, y_test)

By executing this code, we can obtain the  $R^2$  value, which ranges from 0 to 1. A higher  $R^2$  indicates a better fit of the model to the data.

Programming machine learning algorithms using Python allows us to harness the power of AI and ML.  $R^2$  serves as a valuable metric to evaluate the performance of regression models. By understanding the concept of variance and utilizing the scikit-learn library, we can calculate  $R^2$  and assess the accuracy of our models.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: TESTING ASSUMPTIONS**

In this tutorial, we will be discussing the process of testing assumptions in machine learning. Up until this point, we have been focusing on the algorithms and their outputs without thoroughly testing their assumptions. To address this, we will be examining two major algorithms: the equation for the best fit line and the coefficient of determination (r-squared).

In the field of programming, there is a similar concept known as unit testing, where individual components of a program are tested to ensure their functionality. Although our approach will not be exactly the same as unit testing, the idea is similar. We will test various components of our machine learning model using sample data that we can manipulate.

To begin, we will import the 'random' module as we will be using random numbers in our testing. It is worth noting that these numbers are pseudo-random, meaning they are not truly random. However, for the purposes of our testing, they will suffice.

Next, we will define a function called 'create\_dataset' which will take several parameters. The first parameter is the number of data points we want to create. The second parameter is the variance, which determines the variability of the dataset. The third parameter is the step, which specifies the average increase in the Y value per data point. Lastly, we have the correlation parameter, which can be set to 'positive', 'negative', or 'none'.

Within the 'create\_dataset' function, we will first define the objective, which is to return a numpy array of X values and Y values. We will specify the data type as 'float64' to ensure consistency.

To generate the dataset, we will start with an initial value for Y and an empty list to store the Y values. Using a loop, we will iterate through the desired number of data points. Each Y value will be calculated by adding the step value to the previous Y value, with a random value within the specified variance range. The correlation parameter will determine whether the step value is positive, negative, or zero.

Once we have generated the X and Y values, we will return them as numpy arrays with the specified data type.

By visually examining the best fit line and calculating the r-squared value, we can test the assumptions of our machine learning model. If the data appears more linear and the r-squared value is higher, it indicates that our model is performing better. Conversely, if the data is more spread apart and the r-squared value is lower, it suggests that our model may not be accurately representing the relationship.

Testing assumptions is an important step in machine learning. By using sample data and evaluating the best fit line and r-squared value, we can determine the effectiveness of our model. This process allows us to identify any discrepancies and make necessary adjustments to improve the accuracy of our predictions.

To program machine learning algorithms, it is important to test the assumptions made during the process. In this didactic material, we will discuss how to test assumptions using Python.

One assumption we often encounter is the correlation between variables. To generate data with no correlation, we can use the "wise.append" function. By iterating through a range and appending a random value to the current value, we can create data that is somewhat varied but not correlated.

To introduce correlation, we can modify the "wise.append" function. If we want a positive correlation, we can use the "Val += step" statement. On the other hand, if we want a negative correlation, we can use the "Val -= step" statement.

To create a sample dataset, we need both X and Y values. We can use a one-line for loop to generate X values based on the length of the Y values. By calling the "create\_dataset" function and specifying the desired variance, step, and correlation, we can obtain the X and Y values.

To test our assumptions, we can print the coefficient of determination (R-squared) for the dataset. The



coefficient of determination measures how well the data fits the regression line. If the variance is decreased, the coefficient of determination should decrease significantly. Conversely, if the variance is increased, the coefficient of determination should increase.

Additionally, we can change the correlation to test its impact on the dataset. If we set the correlation to false, we should observe an ugly dataset with a low coefficient of determination.

By automating this process, we can write a program that calculates the coefficient of determination for different sample datasets. We can start with a specific variance, change it, and compare the resulting coefficient of determination to the initial value. This approach allows us to test our assumptions and evaluate the impact of different factors on the dataset.

Testing assumptions is an essential step in programming machine learning algorithms. By modifying variables such as variance and correlation, we can observe the effects on the dataset and evaluate the accuracy of our assumptions.

#### Machine Learning with Python - Programming machine learning - Testing assumptions

In machine learning, it is important to test assumptions and ensure that the chosen model is appropriate for the given data set. Linear regression is one commonly used model, but it may not always be suitable for nonlinear data.

If a data set is nonlinear and linear regression is applied, the resulting R-squared value will be very low, indicating that the data cannot be effectively modeled using linear regression. However, there are other forms of classification and machine learning that can be used instead.

When working with large scripts or programs, it is important to ensure that the model is accurate. While visual inspection can help determine the best fit line, R-squared cannot be fully tested in this manner. However, it is possible to create a program that checks whether the R-squared value aligns with the expected assumptions.

Regression is an important concept in machine learning, but there are a few additional points that need to be addressed. Firstly, it is essential to consider a fundamental aspect of machine learning that may be overlooked. Secondly, an error made during the demonstration should be corrected and used as a learning opportunity.

To illustrate these points, let's refer to the code. In the example, the data is modified to represent a 10% shift in price. The blue line represents the prediction line, which includes weekends and holidays, while the stock price data only occurs on weekdays. This modification results in a higher price, but the prediction line remains the same. This demonstrates the impact of linear models and the importance of considering all relevant factors.

Now, let's address the mistakes made during the demonstration. The first mistake was a typo, where a colon was mistakenly added at the end of the variable 'X'. This error does not affect the functionality of the code. The second mistake relates to the slicing of the 'X' variable. Instead of properly defining 'X' as the first 90% of the data for training, it was incorrectly sliced after redefining it. This error results in using only a portion of the intended training data.

To correct this, the proper slicing should be applied to 'X' to ensure that the first 90% is used for training. By making this correction, the model will generate results that closely resemble the previous ones.

Lastly, it is important to discuss the fundamentals of choosing features for training. While the example used in this demonstration focused on simplicity, real-world machine learning problems are often more complex. Therefore, careful consideration should be given to selecting features that directly impact the desired outcome.

When working with machine learning in Python, it is important to test assumptions and select appropriate models for the given data set. Linear regression may not always be suitable for nonlinear data, and other forms of classification and machine learning should be considered. Additionally, it is important to carefully examine and correct any errors in the code, as well as choose features that directly influence the desired outcome.

When programming machine learning algorithms, it is important to test the assumptions made about the data. In this context, the discussion revolves around the impact of different features on the price of a stock. The



transcript highlights that certain features, such as high minus low percent, percent change, and volume, do not directly impact the price. These features are indicators of volatility and magnitude but do not have a direct relationship with the price.

To illustrate this, the transcript suggests dropping the adjusted close feature and observing the effect on the prediction. The result is a flatline, indicating that the features considered are not suitable for predicting stock prices accurately. The transcript then goes on to explain that stock prices are indicative of the overall value of a company. Factors such as quarterly earnings, price to earnings ratio, price to earnings to growth ratio, and book value play a significant role in determining the value of a company. Therefore, to predict stock prices effectively, it is essential to consider features that attempt to predict the company's overall value.

The transcript acknowledges that predicting stock prices can be a complex task due to the numerous factors involved and the dynamic nature of the market. It mentions the availability of a tutorial series for a more in-depth exploration of investing with fundamental features of companies. The series covers concepts such as quarterly earnings, price to earnings to growth ratio, and book value. However, the transcript emphasizes that mistakes are common in programming and encourages learning from them.

When programming machine learning algorithms for stock price prediction, it is important to select features that reflect the overall value of the company. Factors such as quarterly earnings, price to earnings ratio, price to earnings to growth ratio, and book value are more relevant than features like high minus low percent, percent change, and volume. Understanding the fundamental features of companies and their impact on stock prices is essential for accurate predictions.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: INTRODUCTION TO CLASSIFICATION WITH K NEAREST NEIGHBORS**

Classification is an important concept in machine learning, and one of the classification algorithms we will be covering is K nearest neighbors (KNN). In classification, the objective is to create a model that can properly divide or separate our data into different groups.

To better understand this, let's consider an example. Imagine we have a graph with some data points. Visually, we can see that there are two distinct groups of points. This intuitive grouping is known as clustering. However, classification is even simpler than clustering.

In classification, we have a dataset that consists of two groups, let's say pluses and minuses. The goal is to create a model that can accurately fit both groups and properly divide them. Now, suppose we have an unknown data point. Based on its visual proximity to the existing points, we would assign it to either the pluses or minuses group.

For example, if the unknown point is closer to the pluses, we would assign it to the pluses group. Similarly, if it is closer to the minuses, we would assign it to the minuses group. The proximity to the existing points is the key factor in making this decision.

This is where K nearest neighbors comes into play. With KNN, we consider the K number of closest neighbors to the unknown point. For example, if K is equal to 2, we find the two closest points to the unknown point. These two points will then "vote" on the identity of the unknown point. If both points are pluses, the unknown point will be classified as a plus, and vice versa.

It is important to note that the choice of K can affect the classification result. In our example, if K was equal to 2 and the two closest points were one plus and one minus, we would have a split vote. In such cases, additional techniques can be used to break ties and determine the final classification.

KNN is a simple and intuitive classification algorithm. It is particularly useful when dealing with low-dimensional data. However, as the dimensionality increases, it becomes impractical to visually assess the proximity of points. This is where the power of machine learning algorithms comes into play, as they can handle high-dimensional data and make accurate classifications based on proximity.

K nearest neighbors is a classification algorithm that assigns a label to an unknown data point based on its proximity to the K closest points in the dataset. By considering the votes of these neighbors, the algorithm determines the class of the unknown point. The choice of K can impact the classification outcome, and additional techniques may be used to resolve ties.

K nearest neighbors is a simple algorithm used for classification in machine learning. When using this algorithm, it is recommended to choose an odd value for K, such as 3, to avoid ties in the voting process. In the case of  $K=3$ , we would consider the three nearest neighbors to a data point and determine their classes. For example, if two neighbors are negative and one is positive, we would classify the data point as negative.

It is important to note that K nearest neighbors can be applied to datasets with more than two classes. However, if there are three classes, it is not advisable to use  $K=3$ , as it may result in a split vote. In such cases, it is recommended to have at least five neighbors for a more accurate classification.

In addition to providing classification results, K nearest neighbors also allows for assessing the confidence of the classification. For instance, if a data point is classified as negative based on a two out of three vote, we can assign a confidence level of 66% to the classification. This confidence level can be different from the overall accuracy of the trained model.

Despite its simplicity, K nearest neighbors has some limitations. One drawback is the computation of distances between data points, which typically involves calculating the Euclidean distance. This process can be time-consuming, especially for large datasets. Although there are techniques to speed up the computation, the efficiency of the algorithm decreases as the dataset size increases.

Moreover, in K nearest neighbors, there is no clear separation between training and testing phases. The same set of points is used for both tasks. While there are methods to improve the training process, they are beyond the scope of this material.

It is worth mentioning that K nearest neighbors may not scale well for very large datasets, as its performance is inferior to other algorithms, such as support vector machines. However, for classification tasks with datasets up to a gigabyte in size, K nearest neighbors can still provide fast and accurate results. Additionally, the algorithm can be easily parallelized, allowing for efficient computations.

K nearest neighbors is a straightforward algorithm for classification tasks. It considers the K nearest neighbors to a data point and determines its class based on a majority vote. While it has limitations in terms of scalability and training process, it remains a viable option for many classification tasks.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: K NEAREST NEIGHBORS APPLICATION**

In this tutorial, we will continue our discussion on classification in machine learning, specifically focusing on K nearest neighbors (KNN). KNN is a technique that compares the Euclidean distance between a new data point and existing data points to determine its class membership. We will be using scikit-learn for implementing KNN and a dataset from the University of California, Irvine (UCI) for our experiments.

The UCI dataset repository ([archived.ics.uci.edu/ml/datasets](https://archive.ics.uci.edu/ml/datasets)) offers a wide range of datasets categorized based on tasks such as classification, regression, clustering, and more. Each dataset also provides information about attribute types, data types, and the number of instances. For this tutorial, we will be using the breast cancer dataset.

To access the breast cancer dataset, navigate to the UCI website and search for the breast cancer dataset. Once on the dataset page, locate the data folder and download the actual data file. Additionally, you may want to read the accompanying files for more information about the dataset, including citation requests and usage information.

The breast cancer dataset contains various attributes, with the 11th attribute representing the class (benign or malignant). We will consider the class attribute as our target variable for prediction. After downloading the dataset, open the data file and add a header row manually to identify the attributes. Ensure that the class attribute is represented by numerical values (2 for benign and 4 for malignant) instead of strings, as most machine learning algorithms require numerical inputs.

It is worth noting that the dataset contains missing attribute values denoted by a question mark (?). The class distribution indicates that there are 458 benign tumors and 241 malignant tumors, suggesting that the dataset is slightly imbalanced.

Once the dataset is prepared, we can proceed with loading it into our working environment for further analysis and modeling.

To begin programming machine learning with Python, specifically the K nearest neighbors application, we need to handle missing and bad data. First, we import the necessary libraries: numpy (NP) and scikit-learn (SKlearn) for preprocessing, cross-validation, and neighbors. We also import pandas (PD) for data manipulation.

Next, we load the dataset, which is in a text file format, into a dataframe using the `PD.read_csv` function. However, we notice that there are missing data denoted by question marks. To address this, we replace all question marks with a large negative value (-99,999) using the `DF.replace` function. This allows us to treat the missing data as outliers instead of discarding them.

We then identify and remove any useless data from the dataframe. In this case, the ID column does not have any implication on whether a tumor is benign or malignant, so we drop it using the `DF.drop` function.

After preprocessing the data, we need to define our features (X) and labels (Y). The features include all columns except the class column, while the labels are just the class column. We store these in numpy arrays using the `NP.array` function.

To evaluate our model's performance, we perform cross-validation. We split the data into training and testing sets using the `cross_validation.train_test_split` function from scikit-learn. By specifying a test size of 0.2 (20%), 80% of the data is used for training and 20% for testing.

This completes the initial steps of programming machine learning with Python, specifically the K nearest neighbors application. We have handled missing data, removed useless data, defined our features and labels, and performed cross-validation to prepare for training and testing our model.

In this didactic material, we will explore the application of the K nearest neighbors (KNN) algorithm in machine learning using Python. KNN is a classification algorithm that predicts the class of a data point based on the

classes of its nearest neighbors.

To begin, we need to define the classifier. In this case, we will use the `KNeighborsClassifier` class from the `scikit-learn` library. We can create an instance of the classifier by assigning it to a variable, like so:

```
1. classifier = KNeighborsClassifier()
```

Next, we need to fit the classifier to our training data. This can be done using the `fit` method, which takes in the features (`X_train`) and the corresponding labels (`y_train`). The code for fitting the classifier would look like this:

```
1. classifier.fit(X_train, y_train)
```

Once the classifier is fitted, we can evaluate its performance by testing it on the test data. This can be done using the `score` method, which calculates the accuracy of the classifier on the test data. The code for calculating the accuracy would look like this:

```
1. accuracy = classifier.score(X_test, y_test)
```

The accuracy is a measure of how well the classifier performs, with higher values indicating better performance. In this case, we are using the term "accuracy" instead of "confidence" to measure the performance of the classifier. It is important to note that KNN also has a concept of confidence, which is related to the voting mechanism used in the algorithm.

In our example, the accuracy of the classifier is 96.4%. This means that the classifier correctly predicted the class of 96.4% of the test data points. However, depending on the application, higher accuracy may be desired. For example, in a self-driving car scenario, it is important to have a high accuracy to differentiate between objects on the road.

We can further improve the accuracy of the classifier by considering additional features. In this case, including the ID column significantly improves the accuracy of the classifier. This highlights the importance of feature selection and engineering in machine learning.

Finally, once the classifier is trained, we can use it to make predictions on new, unseen data points. To do this, we can use the `predict` method, which takes in the features of the new data point and returns the predicted class. The code for making a prediction would look like this:

```
1. prediction = classifier.predict(example_measures)
```

In our example, we created a new example with some arbitrary feature values and made a prediction using the trained classifier. The predicted class is then printed to the console.

It is worth mentioning that in some cases, it is common to save the trained classifier to a file using the `pickle` module. This allows for reusing the classifier without having to train it again. However, in this case, the classifier is trained and tested quickly, so saving it is not necessary.

We have explored the application of the K nearest neighbors algorithm in machine learning using Python. We learned how to define and fit the classifier, evaluate its performance using accuracy, and make predictions on new data points. Feature selection and engineering play an important role in improving the accuracy of the classifier.

In this didactic material, we will discuss the application of the K nearest neighbors algorithm in machine learning using Python. We will focus on working with `scikit-learn` and `numpy` libraries, specifically in reshaping and handling different shapes of data.

To begin, let's consider a scenario where we have one sample with a value of -1. If we were to have two samples, we would need to represent them as a list of lists. In this case, we can create a list of lists by copying the existing sample and changing the value to 2. By doing this, we now have two samples.

Now, let's address the situation where we have an unknown number of predictions. For example, one week we may have 15 patients, and the next week we may have 45 patients. To handle this variability, we need to dynamically reshape the data without hard coding the shape every time.

To achieve this, we can convert the list of lists to a numpy array. Once we have the numpy array, we can use the reshape function. Instead of specifying a fixed number, we can replace it with the length of the example measures. This allows us to automatically reshape the data to match the desired shape that scikit-learn expects.

By following this approach, we can programmatically predict on any number of predictions without the need for manual adjustments. It is important to note that this technique is not specific to K nearest neighbors but can be applied to any scikit-learn classifier.

We have discussed the process of reshaping data using scikit-learn and numpy libraries. By converting the data to a numpy array and using the reshape function, we can dynamically handle different shapes of data, allowing us to feed it into scikit-learn classifiers seamlessly.

In real-world examples, we have observed that the K nearest neighbors algorithm can achieve prediction accuracies ranging from 94% to 98%. This level of accuracy is impressive, considering the simplicity of the algorithm. In the next video, we will take a step further and write our very own K nearest neighbors algorithm from scratch. We will then compare its performance to scikit-learn.

Please feel free to leave any questions or comments below. Thank you for watching and for your continued support.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: EUCLIDEAN DISTANCE**

In this tutorial, we will be discussing the concept of Euclidean distance in the context of machine learning with Python. Euclidean distance is a fundamental concept that forms the basis of the K nearest neighbors algorithm, which we have been exploring in the previous tutorials.

Euclidean distance is named after Euclid, a famous mathematician known as the father of geometry. It is defined as the square root of the sum of the squared differences between corresponding coordinates of two points. In other words, it measures the straight-line distance between two points in a multidimensional space.

To calculate the Euclidean distance between two points, we follow a simple formula. Let's say we have two points, Q and P, with coordinates (Q1, Q2, ..., Qn) and (P1, P2, ..., Pn) respectively. The Euclidean distance between these two points can be calculated as follows:

$$\text{distance} = \sqrt{(Q_1 - P_1)^2 + (Q_2 - P_2)^2 + \dots + (Q_n - P_n)^2}$$

Here, n represents the number of dimensions in the data. The formula involves taking the difference between each corresponding coordinate, squaring it, summing up all the squared differences, and finally taking the square root of the sum.

To illustrate this with an example, let's consider two points: Q(1, 3) and P(2, 5). These points have two dimensions. Therefore, the Euclidean distance between them can be calculated as:

$$\text{distance} = \sqrt{(1 - 2)^2 + (3 - 5)^2}$$

$$\text{distance} = \sqrt{(-1)^2 + (-2)^2}$$

$$\text{distance} = \sqrt{1 + 4}$$

$$\text{distance} = \sqrt{5}$$

$$\text{distance} \approx 2.236$$

Now, let's implement the calculation of Euclidean distance in Python. We can use the math module and the sqrt function to calculate the square root. Here is an example code snippet:

1.	from math import sqrt
2.	
3.	plot1 = (1, 3)
4.	plot2 = (2, 5)
5.	
6.	distance = sqrt((plot1[0] - plot2[0])**2 + (plot1[1] - plot2[1])**2)
7.	
8.	print("The Euclidean distance between plot1 and plot2 is:", distance)

In the code above, we import the sqrt function from the math module. We define the coordinates of the two points as tuples, plot1 and plot2. Then, we calculate the Euclidean distance using the formula discussed earlier. Finally, we print the result.

By understanding and implementing the Euclidean distance calculation, we gain an important step in building the foundation for various machine learning algorithms, such as K nearest neighbors.

The Euclidean distance is a fundamental concept in machine learning, specifically in the K nearest neighbors algorithm. It allows us to measure the distance between two points in a multi-dimensional space. In this didactic material, we will learn how to calculate the Euclidean distance using Python.

To calculate the Euclidean distance, we first need to have two points in a multi-dimensional space. Each point is represented by a set of coordinates. The Euclidean distance between these two points is the square root of the sum of the squared differences of their corresponding coordinates.

Let's consider an example. Suppose we have two points, A and B, in a two-dimensional space. Point A has coordinates (x1, y1), and point B has coordinates (x2, y2). The Euclidean distance between A and B can be calculated using the following formula:

$$\text{distance} = \sqrt{(x2 - x1)^2 + (y2 - y1)^2}$$

To demonstrate this calculation in Python, we can define a function that takes the coordinates of two points as input and returns the Euclidean distance. Here is an example implementation:

1.	import math
2.	
3.	def euclidean_distance(x1, y1, x2, y2):
4.	distance = math.sqrt((x2 - x1) <sup>2</sup> + (y2 - y1) <sup>2</sup> )
5.	return distance
6.	
7.	# Example usage
8.	point_a = (1, 2)
9.	point_b = (4, 6)
10.	distance = euclidean_distance(point_a[0], point_a[1], point_b[0], point_b[1])
11.	print(distance)

In this example, we calculate the Euclidean distance between point A with coordinates (1, 2) and point B with coordinates (4, 6). The result will be printed as 5.0.

By using this formula and the provided Python code, you can calculate the Euclidean distance between any two points in a multi-dimensional space. This distance metric is widely used in machine learning algorithms, such as K nearest neighbors, to measure the similarity between data points.

In the next tutorial, we will focus on creating the framework for the K nearest neighbors algorithm. This framework will take a dataset and use the Euclidean distance to classify points. Stay tuned for the upcoming tutorial!



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: DEFINING K NEAREST NEIGHBORS ALGORITHM**

In this tutorial, we will be discussing the K nearest neighbors algorithm and how to implement it in Python for machine learning. The K nearest neighbors algorithm is a simple yet powerful classification algorithm that predicts the class of a given point based on the vote of its K nearest neighbors.

To begin, we need to import the necessary libraries. We will be using numpy for its efficient mathematical operations, matplotlib for data visualization, and collections for vote counting. We will also import warnings to handle invalid values for K.

Next, we will define a dataset consisting of two classes, K and R. Each class has a set of features, which are two-dimensional coordinates. We will represent the features as a list of lists. For example, the class K has features (1, 2), (2, 3), and (3, 1), while the class R has features (6, 5), (7, 7), and (8, 6).

Now, let's consider a new point with features (5, 7). We want to determine which class this point belongs to. Visually, it seems like it belongs to class R, but we will write an algorithm to make this determination.

To visualize the dataset, we will use a scatter plot. We will iterate through each class and its corresponding features, and scatter plot each feature with a different color. We will use the scatter function from matplotlib to accomplish this.

Finally, we will display the scatter plot using the show function from matplotlib. The scatter plot will show the two classes, K and R, with different colors.

In the next tutorial, we will write the K nearest neighbors algorithm and apply it to a more realistic dataset. This simple dataset was used for illustrative purposes only.

The K nearest neighbors (KNN) algorithm is a simple and intuitive machine learning algorithm used for classification and regression tasks. In this tutorial, we will define the KNN algorithm and start building it step by step.

Before we dive into the algorithm, let's visualize our data. We have scattered the data using the PLT dot scatter function. By looking at the scatter plot, we can visually identify the groups to which the data points belong. In this case, we can see that the data points belong to the red group.

Now, let's move on to defining the KNN algorithm. To calculate the K nearest neighbors, we need to pass through the training data, the data we are trying to predict, and a value for K. We will create a function for this purpose.

In the function, we will check if the length of the data is greater than or equal to the chosen value of K. If it is, we will send a warning to the user. This is because having a value of K greater than the total number of voting groups may lead to inaccurate predictions.

Next, we will create a variable called "vote\_results" to store the results of the KNN algorithm. Finally, we will return the vote\_results.

Please note that this is just the starting point of our KNN algorithm. In the next tutorial, we will continue building it to the point where we can pass in the data and obtain the vote result. We will also test it on real data.

If you have any questions or concerns, please feel free to post them below. Thank you for watching and for your support. Stay tuned for the next tutorial!

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: PROGRAMMING OWN K NEAREST NEIGHBORS ALGORITHM**

In this part of our machine learning tutorial series, we will focus on programming our own K nearest neighbors algorithm using Python. So far, we have created a function that warns the user when they attempt to perform an invalid action. Now, we need to move on to the actual implementation of the K nearest neighbors algorithm.

The first step in this process is to calculate the K nearest neighbors. To do this, we need to compare a given data point to all other data points in the dataset. This is the main challenge of the K nearest neighbors algorithm. One approach to address this challenge is to use a radius, which allows us to look within a certain distance of a point and ignore outliers. By doing this, we can simplify the calculation of the Euclidean distance, which is used to determine the proximity between data points.

To implement the K nearest neighbors algorithm, we start by creating a list of lists to store the distances between data points. We iterate through each group or class in the dataset and then iterate through the features of each data point within that group. For each pair of data points, we calculate the Euclidean distance using the formula: square root of  $((\text{feature } 0 - \text{predict } 0)^2 + (\text{feature } 1 - \text{predict } 1)^2)$ . This formula compares the features of the prediction point with the features of each data point.

However, this approach is not very efficient and is limited to two-dimensional data. To address this, we can use the numpy library to perform faster calculations and handle data with any number of dimensions. By using numpy functions such as square root, sum, and arrays, we can rewrite the Euclidean distance calculation as a more concise and efficient expression.

Here is an example of the numpy-based calculation of Euclidean distance: `Euclidean distance = numpy.linalg.norm(features - predict)`. This expression uses the `numpy.linalg.norm` function to calculate the Euclidean distance between the features of a data point and the prediction point. This approach is faster and more flexible than the previous method.

Finally, we append the calculated Euclidean distance and the corresponding group to the distances list. This allows us to keep track of the distances between data points and their respective groups.

To summarize, in this part of the tutorial, we have discussed the implementation of the K nearest neighbors algorithm. We have explored the challenges of comparing data points and calculating the Euclidean distance. We have also demonstrated two different ways to calculate the Euclidean distance, one using basic Python operations and the other using the numpy library for faster and more flexible calculations.

In this didactic material, we will discuss the topic of programming our own K nearest neighbors algorithm for machine learning using Python. This algorithm is a popular and simple classification algorithm that can be used to predict the class of a new data point based on its proximity to existing data points.

To begin, let's first understand the concept of distance. In our algorithm, distance refers to the measure of similarity or dissimilarity between two data points. We will represent distance as a list of lists, where the first item in the list is the actual distance and the second item is the group to which the data point belongs.

To implement the algorithm, we will use a one-liner for loop. We will sort the distances and select the top K distances. Once we have the top K distances, we will only consider the group to which each distance belongs. This will allow us to rank the distances and determine the most common group among the top K distances.

Next, we will use the Counter function from the collections module to count the occurrences of each group in the list of votes. We will retrieve the most common group, which is the group with the highest count, and return it as the result.

To test our algorithm, we can pass a dataset and new features to predict on. We will set the value of K to determine the number of nearest neighbors to consider. The algorithm will return the most voted group for each data point.

It is important to note that our algorithm is a simplified version of the K nearest neighbors algorithm and may not perform as well as the implementation provided by scikit-learn, a popular machine learning library in Python. However, it serves as a good starting point for understanding the inner workings of the algorithm.

We have successfully programmed our own K nearest neighbors algorithm using Python. We have discussed the concept of distance, implemented the algorithm using a one-liner for loop, and tested it on a dataset. While our implementation may not be as robust as the one provided by scikit-learn, it provides a solid foundation for understanding the algorithm.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: APPLYING OWN K NEAREST NEIGHBORS ALGORITHM**

In this tutorial, we will be applying our own K nearest neighbors algorithm to a real-world dataset. Specifically, we will use the breast cancer dataset and compare our accuracy to that of the scikit-learn library. The goal is to determine if our algorithm performs similarly or if scikit-learn's classifier outperforms ours.

To start, we need to clean up some code. We will remove unnecessary information and matplotlib graphs since the dataset has too many dimensions to be graphed effectively. Additionally, we will import the pandas library as PD and the random library to load and shuffle the dataset.

Next, we will read the breast cancer dataset using the PD.read\_csv() function. We will replace any question marks in the dataset with a value of -99999. This is important because K nearest neighbors relies on distance calculations, and missing data can significantly impact the results. We will also drop the ID column, as it does not provide any useful information.

After cleaning the dataset, we will convert it to a list of lists using the .values.tolist() function. This step is necessary because some values in the dataset were being treated as strings instead of integers or floats. By converting everything to floats, we ensure consistency and compatibility with our algorithm.

Now that we have the cleaned and converted dataset, we can shuffle it using the random.shuffle() function. Shuffling the dataset is possible because we have converted it to a list of lists, preserving the relationship between features and labels. We will print a subset of the shuffled data to demonstrate the shuffling process.

In this tutorial, we applied our own K nearest neighbors algorithm to the breast cancer dataset. We cleaned the dataset, removed unnecessary information, and converted it to a list of lists. We then shuffled the dataset and printed a subset of the shuffled data. This allows us to compare the accuracy of our algorithm to that of the scikit-learn library.

In machine learning, one common task is to apply the K nearest neighbors algorithm. This algorithm is used for classification tasks, where we want to predict the class of a new data point based on its neighboring data points. In this didactic material, we will go through the process of programming our own K nearest neighbors algorithm using Python.

First, we need to prepare our data. We start by shuffling our dataset, which ensures that our training and test sets are representative of the entire dataset. To split our data into training and test sets, we use a test size of 0.2, meaning that 20% of the data will be used for testing. We create empty lists for our train and test sets, and then slice our shuffled data accordingly.

Next, we need to populate dictionaries for our train and test sets. We iterate through the train data and append the elements to the train set dictionary. The last element in each list represents the class, either 2 or 4, where 2 represents benign and 4 represents malignant. We use negative indexing to access the last element and append the rest of the elements up to the last element.

We repeat the same process for the test data, populating the test set dictionary.

Now, we are ready to pass the information to our K nearest neighbors algorithm. We create counters for correct and total predictions. We iterate through each group in the test set and then iterate through the data in each group. We pass the data and the train set dictionary to our K nearest neighbors algorithm, with a value of K equal to 5.

To determine if our predictions are correct, we compare the predicted group with the actual group from the test set. If they are equal, we increment the correct counter. In either case, we increment the total counter.

Finally, we calculate the accuracy by dividing the correct counter by the total counter. We print the accuracy to evaluate the performance of our own K nearest neighbors algorithm.

---

**EUROPEAN IT CERTIFICATION CURRICULUM SELF-LEARNING PREPARATORY MATERIALS**

---

In this tutorial, we have implemented our own K nearest neighbors algorithm and applied it to a dataset. We achieved an accuracy of 97.8% on our first run and 95.6% on the second run.

Now, let's compare our results with scikit-learn, a popular machine learning library in Python. By doing this, we can evaluate the performance of our algorithm against a well-established and widely-used implementation.

In the next tutorial, we will also calculate the confidence of our model. This will provide us with a measure of how confident we can be in the predictions made by our algorithm.

If you have any questions, comments, or concerns up to this point, please feel free to leave them below. Otherwise, in the next tutorial, we will compare our algorithm with scikit-learn and calculate the confidence of our model.

Thank you for watching and for all your subscriptions. Until next time!

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: PROGRAMMING MACHINE LEARNING****TOPIC: SUMMARY OF K NEAREST NEIGHBORS ALGORITHM**

K nearest neighbors is a machine learning algorithm used for classification tasks. In this algorithm, the value of K represents the number of nearest neighbors used to make predictions. In this tutorial, we will discuss the concept of K accuracy and predictions.

One question that arises is whether increasing the value of K would necessarily lead to an increase in accuracy. To explore this, let's consider an example where K is set to 25. Running the algorithm multiple times with this value of K, we observe that the accuracy fluctuates between 88% and 98%. This indicates that increasing K does not always guarantee higher accuracy. Similarly, when we set K to 75, the accuracy remains consistently high at around 95-97%.

It is worth noting that the dataset used in this example has about 600 data points, with only 30% being malignant and the remaining 70% being benign. This skewed distribution may affect the accuracy of the algorithm. Increasing K to 200, which includes more data points, actually results in lower accuracy. This is because the algorithm is influenced by the majority class (benign) due to the imbalanced distribution.

Hence, it is important to carefully choose the value of K based on the dataset and the problem at hand. While 5 appears to be a good guess in this case, it is recommended to experiment with different values of K to determine the optimal choice for a given dataset.

Moving on, let's discuss the concept of confidence versus accuracy in K nearest neighbors. Accuracy measures whether the classification is correct, while confidence provides information about the certainty of the classification. The classifier can assign a confidence score based on the votes obtained from the nearest neighbors.

To calculate confidence, we divide the number of votes for a class by the value of K. For example, if K is set to 5, we hope to have a confidence score of 5. By printing the confidence scores along with the classification results, we can gain insights into the confidence level of the predictions.

In the provided code, the confidence scores are displayed for each classification result. Most of the correct predictions have a confidence score of 1.0, indicating high certainty. However, there are some incorrect predictions with lower confidence scores, such as 0.8 and 0.6.

By adjusting the test size, we can observe changes in the confidence scores. Decreasing the test size sacrifices more data, potentially leading to a decrease in accuracy. However, it also reduces the number of instances with 100% confidence, highlighting the importance of considering confidence levels in decision-making.

K nearest neighbors is a versatile algorithm for classification tasks. The choice of K can significantly impact the accuracy of the algorithm, and it is important to consider the distribution of classes in the dataset. Additionally, confidence scores provide valuable insights into the certainty of the predictions, allowing for informed decision-making.

The K nearest neighbors algorithm is a popular machine learning algorithm used for classification tasks. In this summary, we will discuss the implementation of the K nearest neighbors algorithm using Python.

To begin, let's review the results obtained from running the algorithm. We performed 10 tests, with 5 tests being run initially. The average accuracy achieved from these tests was 96.2%. It is important to note that the value of k used in these tests was 5.

Next, we made some modifications to the code to ensure accurate results. We removed unnecessary print statements and ensured that the value of k remained consistent. After running the modified code 25 times, we obtained an average accuracy of 96.4%.

Moving on, we explored another version of the K nearest neighbors algorithm. We applied the same modifications to this version and ran it 25 times. The average accuracy achieved in this case was 96.8%.

Now, let's discuss the differences between the two versions. Firstly, the K nearest neighbors algorithm has a default parameter called "n\_jobs", which determines the number of parallel jobs to run. By default, this value is set to 1. However, it can be set to -1 to utilize as many parallel jobs as possible, thereby improving performance.

Additionally, the K nearest neighbors algorithm also has a parameter called "radius", which allows for the exclusion of points outside a specified radius. This can be useful in certain scenarios.

It is worth mentioning that the accuracy achieved by our implementation was comparable to that of other libraries, such as scikit-learn. While our tutorial focuses on machine learning rather than high-performance computing, it is important to note that the K nearest neighbors algorithm can scale well to large datasets.

Furthermore, the K nearest neighbors algorithm is versatile and can be used on both linear and nonlinear data. This flexibility makes it a valuable tool for classification tasks.

The K nearest neighbors algorithm is an effective machine learning algorithm for classification tasks. It can be implemented using Python and provides accurate results. By adjusting parameters such as "n\_jobs" and "radius", the algorithm's performance can be further optimized. Additionally, it is important to consider the nature of the data being analyzed, as the algorithm can handle both linear and nonlinear data.

The K nearest neighbors algorithm is a classification algorithm that can also be used for regression. It works by finding the K closest labeled data points in the training set to a new input and then classifying or predicting the value of that input based on the majority vote or average of the K neighbors.

In the case of classification, if we have an unknown data point, we can measure the squared error between the regression lines of the different classes. The class with the lesser squared error would be assigned to the new data point. This method can be used for linear data, allowing for classification even when forecasting is not required.

However, for datasets with nonlinear data, using regression lines to classify would result in poor accuracy. In such cases, the K nearest neighbors algorithm can be applied. It involves measuring the distance between the new data point and its K closest neighbors. By taking a majority vote among these neighbors, the algorithm can classify the new data point.

For example, consider a dataset with orange and blue dots. Even if there is a best fit line for both classes, the coefficient of determination would be low, indicating poor confidence in the algorithm's classification. However, by using K nearest neighbors, we can classify a new data point by measuring the distance between it and its closest neighbors. If K is equal to three, we would consider the three closest points and take a vote. In this case, if two out of the three closest points are orange, we would classify the new data point as orange.

The K nearest neighbors algorithm has the advantage of being able to handle nonlinear data, making it suitable for classification tasks in such cases. It is a versatile and widely used algorithm in machine learning.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SUPPORT VECTOR MACHINE INTRODUCTION AND APPLICATION**

A support vector machine (SVM) is a supervised machine learning algorithm that is used for classification tasks. It was created by Lotfi A. Zadeh in the 1960s but gained popularity in the 1990s when it was shown to outperform neural networks in tasks such as handwritten number recognition.

The SVM works in vector space and is a binary classifier, meaning it separates data into two groups at a time. It aims to find the best separating hyperplane, also known as the decision boundary, that separates the positive and negative groups. The positive and negative groups can be thought of as two different classes or categories.

To illustrate this, let's consider an example with two groups of data: positive and negative. The objective of the SVM is to find the best hyperplane that maximizes the distance between the hyperplane and the closest data points from both groups. This distance is known as the margin.

In a two-dimensional space, the hyperplane would appear as a line. However, it is important to note that SVMs can work in higher-dimensional spaces as well. The best hyperplane is determined by calculating the perpendicular distance between the hyperplane and the closest data points from each group. The hyperplane with the largest margin is considered the best separating hyperplane.

Once the best separating hyperplane is acquired, the SVM can classify unknown data points. If an unknown data point falls on the right side of the hyperplane, it is classified as a positive sample. Conversely, if it falls on the left side, it is classified as a negative sample.

It is worth mentioning that SVMs are not limited to linear separation. They can also handle non-linear separation by using techniques such as kernel functions.

A support vector machine is a powerful machine learning algorithm used for classification tasks. It separates data into two groups by finding the best separating hyperplane that maximizes the margin between the hyperplane and the closest data points from each group. SVMs can handle both linear and non-linear separation and are widely used in various applications.



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: UNDERSTANDING VECTORS**

In this didactic material, we will discuss the concept of vectors in the context of support vector machines in machine learning. Vectors are an essential component of support vector machines and understanding them is important for building and applying these models effectively.

A vector represents a point in a vector space or feature space. It consists of multiple dimensions, denoted as  $x_1$ ,  $x_2$ , and so on. These dimensions can also be interpreted as features or variables. For example, in a two-dimensional vector space,  $x_1$  and  $x_2$  represent the two dimensions.

To visualize a vector, we can draw a coordinate system with ticks representing the values of each dimension. For instance, if we have a five-unit span in each dimension, the ticks would be labeled as 1, 2, 3, 4, and 5. By plotting the values of a vector on this coordinate system, we can determine its direction and magnitude.

A vector is denoted by a letter, such as vector  $a$ , and is represented by an arrow above the letter. However, sometimes the arrow is omitted, or a bar is placed above the letter to indicate that it is a vector. In our case, we will use the arrow notation.

A vector has both magnitude and direction. The magnitude of a vector is the length of the vector and is denoted by double bars around the vector. The magnitude is calculated using the formula: square root of  $(x_1^2 + x_2^2 + \dots)$ . For example, if we have a vector  $a$  with values 3 and 4 for  $x_1$  and  $x_2$  respectively, the magnitude of vector  $a$  is calculated as the square root of  $(3^2 + 4^2)$ , which equals 5.

The magnitude of a vector can be visualized using the Pythagorean theorem. By treating the vector as the hypotenuse of a right-angled triangle, the lengths of the vector's dimensions can be used as the lengths of the triangle's sides. Applying the Pythagorean theorem, we can calculate the magnitude of the vector.

In addition to magnitude, vectors can be multiplied using the dot product operation. The dot product is the sum of the products of corresponding elements in two vectors. For example, if we have vector  $a$  with values 1 and 3, and vector  $b$  with values 2 and 4, the dot product of  $a$  and  $b$  is calculated as  $(1*2 + 3*4)$ , which equals 10. The dot product results in a scalar value.

Understanding vectors and their properties, such as magnitude and dot product, is important for comprehending the support vector machine algorithm. Support vector machines utilize vectors to represent data points and make predictions based on their relationships. The magnitude of a vector provides information about the length or magnitude of the data point, while the dot product allows us to measure the similarity or dissimilarity between vectors.

In the next part of this educational material, we will delve deeper into the support vector machine algorithm and explore how vectors and dot products are utilized in this context.

Support Vector Machines (SVM) are a powerful machine learning algorithm used for classification and regression tasks. In this material, we will focus on understanding vectors in the context of SVM.

In SVM, vectors play an important role in representing data points. Each data point is represented as a vector in a high-dimensional space. These vectors are used to classify data points into different classes. The goal of SVM is to find an optimal hyperplane that separates the data points of different classes with the maximum margin.

To understand vectors in SVM, it is important to grasp the concept of feature space. In the feature space, each dimension represents a feature or attribute of the data. For example, in a dataset of images, each dimension could represent the intensity of a specific pixel.

Vectors in SVM are used to represent data points in the feature space. Each vector consists of values corresponding to the features of a data point. The number of dimensions in the vector is equal to the number of features in the dataset. These vectors can be represented as points in the feature space.

The hyperplane in SVM is defined by a vector called the normal vector. This vector is orthogonal to the hyperplane and determines its orientation. The normal vector is calculated using the support vectors, which are a subset of the training data points that lie closest to the hyperplane.

Support vectors are the key elements in SVM. They are the data points that define the decision boundary between different classes. The distance between the support vectors and the hyperplane is known as the margin. The goal of SVM is to find the hyperplane with the maximum margin, which provides the best separation between classes.

Vectors in SVM represent data points in the feature space. They are used to calculate the hyperplane and determine the decision boundary between classes. Support vectors are the data points that define the margin and play an important role in finding the optimal hyperplane.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SUPPORT VECTOR ASSERTION**

Support Vector Machines (SVM) are a powerful machine learning algorithm used for classification tasks. In SVM, a decision boundary, known as a separating hyperplane, is created to divide the data points into different classes. But how does SVM actually classify new points after being trained?

To understand this, let's visualize a vector space with positive and negative data points. The goal of SVM is to create a decision boundary that separates these points. Once the decision boundary is established, SVM classifies new points by projecting them onto a vector perpendicular to the separating hyperplane.

Let's consider an unknown data point represented as vector  $u$ . We project vector  $u$  onto the perpendicular vector, denoted as vector  $W$ . By analyzing the position of the projection, we can determine which side of the hyperplane the unknown point falls on. If the projection is on the left side (or any specified side), the point is classified as positive. Conversely, if the projection is on the right side (or any specified side), the point is classified as negative.

To formalize this classification process, we use the equation:  $\text{vector } u \cdot \text{vector } W + B$ , where  $B$  represents the bias term. If the result of this equation is greater than or equal to zero, the point is classified as positive. If the result is less than zero, the point is classified as negative.

However, what if the result of the equation is exactly zero? In this case, the point lies on the decision boundary. So, in summary, SVM classifies new points based on the result of the equation  $\text{vector } u \cdot \text{vector } W + B$ . If the result is greater than or equal to zero, it is classified as positive. If the result is less than zero, it is classified as negative. And if the result is exactly zero, it lies on the decision boundary.

To find the values of vector  $W$  and  $B$ , we need to solve the equation  $\text{vector } u \cdot \text{vector } W + B$ . We already know the values of vector  $u$  (unknown data point), but we need to determine the values of vector  $W$  and  $B$ . These values come with certain constraints, which we will discuss next.

The constraints are derived from the equation  $X - \text{support vector} \cdot \text{vector } W + B = -1$  and  $X + \text{support vector} \cdot \text{vector } W + B = 1$ , where  $X$  represents the feature set and support vector represents the vectors that lie on the decision boundary. We introduce  $Y_{\text{sub } I}$ , which represents the class of the features. If the class is positive,  $Y_{\text{sub } I}$  equals 1. If the class is negative,  $Y_{\text{sub } I}$  equals -1.

By incorporating these constraints, we can formulate an equation to locate the support vectors and solve for vector  $W$  and  $B$ . This equation allows us to accurately classify new points based on their position relative to the decision boundary.

Support Vector Machines use a decision boundary to classify new points. By projecting new points onto a perpendicular vector, SVM determines their position relative to the decision boundary. The classification is based on the result of the equation  $\text{vector } u \cdot \text{vector } W + B$ . If the result is greater than or equal to zero, the point is positive. If the result is less than zero, the point is negative. And if the result is exactly zero, it lies on the decision boundary.

In support vector machines, we use equations to identify the positive and negative support vectors. For the plus class, the equation is  $X_{\text{sub } I} \cdot W + B = 1$ . And for the minus class, it is  $X_{\text{sub } I} \cdot W + B = -1$ .

To simplify the equations, we can represent  $Y_{\text{sub } I}$  as one for the plus class and negative one for the minus class. Multiplying both sides of the equations by  $Y_{\text{sub } I}$ , we get  $Y_{\text{sub } I} \cdot X_{\text{sub } I} \cdot W + B = Y_{\text{sub } I}$ . Since  $Y_{\text{sub } I}$  is either one or negative one, when we multiply, we still get one on both sides of the equation.

Now, let's set both equations equal to zero. To do this, we subtract one from both sides of the top equation and add negative one to the right side of the equation. This results in  $Y_{\text{sub } I} \cdot X_{\text{sub } I} \cdot W + B - Y_{\text{sub } I} = 0$ .

Similarly, for the bottom equation, we subtract one from both sides, resulting in  $Y_{sub\ I} \text{ multiplied by } X_{sub\ I} W \text{ plus } B \text{ minus one equals zero.}$

Therefore, the equation to derive support vectors for both the positive and negative classes is  $Y_{sub\ I} \text{ multiplied by } X_{sub\ I} W \text{ plus } B \text{ minus one equals zero.}$

In the next tutorial, we will discuss what we can do with support vectors once we have identified them. If you have any questions or comments, please feel free to leave them below. Thank you for watching and for your support and subscriptions.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SUPPORT VECTOR MACHINE FUNDAMENTALS**

Support Vector Machine (SVM) is a powerful machine learning algorithm used for classification and regression tasks. In this tutorial, we will focus on the fundamentals of SVM and how it works.

Support vectors are the key components of SVM. These are the data points that lie closest to the decision boundary, also known as the hyperplane, and play an important role in determining the optimal hyperplane for classification. The goal of SVM is to find the hyperplane that maximizes the margin, which is the distance between the support vectors on either side of the decision boundary.

To calculate the width of the margin, we use the equation:  $\text{width} = (x_+ - x_-) \cdot w / \|w\|$ , where  $x_+$  and  $x_-$  are the support vectors,  $w$  is the weight vector, and  $\|w\|$  is the magnitude of  $w$ . The objective is to maximize this width.

To simplify the equation, we can express the width as:  $\text{width} = 2 / \|w\|$ . This means that to maximize the width, we need to minimize the magnitude of the weight vector  $w$ .

To further simplify the problem, we can minimize half the magnitude of vector  $w$  squared, which is equivalent to minimizing the magnitude of vector  $w$ . This mathematical convenience allows us to plug the equation into the Lagrangian.

The Lagrangian is a function that incorporates the constraints of the problem. In the case of SVM, the constraint is defined as:  $y \cdot (x \cdot w + b) - 1 = 0$ , where  $y$  is the class label,  $x$  is the feature vector,  $w$  is the weight vector, and  $b$  is the bias term.

By introducing Lagrange multipliers, we can solve the optimization problem and find the optimal values for  $w$  and  $b$  that minimize the magnitude of vector  $w$  while satisfying the constraint equation.

The support vector machine algorithm aims to find the optimal hyperplane that maximizes the margin between the support vectors. This is achieved by minimizing the magnitude of the weight vector  $w$ , subject to the constraint equation involving the Lagrange multipliers.

By understanding the fundamentals of support vector machines, we can apply this powerful algorithm to various classification and regression problems.

Support Vector Machines (SVM) are a powerful machine learning algorithm used for classification and regression tasks. In this didactic material, we will focus on the fundamentals of SVM and its optimization process.

The main goal of SVM is to find the best hyperplane that separates the data points of different classes with the maximum margin. To achieve this, SVM uses a mathematical formulation called the Lagrangian, which involves Lagrange multipliers. The Lagrangian consists of two parts: the first part is the magnitude of the weight vector, denoted as  $\|w\|$ , and the second part is the sum over the Lagrange multipliers, denoted as  $\sum \alpha_i$ .

The constraint in SVM is defined by the support vectors, which are the data points closest to the decision boundary. The equation for the hyperplane in SVM is given by  $[w] \cdot [x] + b = 0$ , where  $[w]$  is the weight vector,  $[x]$  is the input vector, and  $b$  is the bias term. By modifying the bias term, we can shift the hyperplane up or down.

To optimize the SVM, we need to differentiate the Lagrangian with respect to the weight vector  $[w]$  and the bias term  $b$ . Differentiating the Lagrangian with respect to  $[w]$  gives us the equation  $[w] = \sum (\alpha_i \cdot y_i \cdot x_i)$ , where  $\alpha_i$  is the Lagrange multiplier associated with the  $i$ -th data point,  $y_i$  is the corresponding class label, and  $x_i$  is the input vector.

Differentiating the Lagrangian with respect to  $b$  gives us the equation  $\sum (\alpha_i \cdot y_i) = 0$ . These equations form the basis for solving the SVM optimization problem.

The optimization problem in SVM is a quadratic programming problem, which involves maximizing the

Lagrangian with respect to the Lagrange multipliers  $\alpha_i$ . The Lagrange multipliers play an important role in determining the support vectors and the decision boundary.

One of the downsides of SVM is its complexity, both in terms of mathematics and the optimization problem itself. Another drawback is that all the feature sets need to be in memory for optimization, which may not be feasible for large datasets. However, there are methods like sequential minimal optimization (SMO) that can be used to work with smaller or larger datasets.

SVM is a powerful machine learning algorithm that uses the concept of hyperplanes to separate data points of different classes. It involves the optimization of the Lagrangian using Lagrange multipliers. Despite its complexity and memory requirements, SVM has proven to be effective in various applications.

Support Vector Machines (SVMs) are widely used in the field of Artificial Intelligence and Machine Learning. They are known for their ability to handle complex datasets and make accurate predictions. In this didactic material, we will discuss the fundamentals of Support Vector Machines using Python.

Once a Support Vector Machine is trained, it can be used to classify new data points. The classification process involves multiplying the weight vector ( $W$ ) with the input vector ( $x$ ) and adding a bias term ( $B$ ). The result of this calculation is then passed through a sign function. If the result is positive, the data point belongs to the positive class, and if it is negative, the data point belongs to the negative class.

One of the advantages of SVMs is that once the machine is trained, the original features used for training are no longer needed. This makes SVMs efficient in terms of memory usage. SVMs can handle high-dimensional data and are effective in separating data points using hyperplanes in the feature space.

In the next material, we will dive deeper into the concepts of SVMs and their constraints. We will also simplify the problem to enhance understanding. We will then proceed to write our own Support Vector Machine from scratch using Python. This will involve creating an optimization algorithm to train the SVM and perform predictions.

By implementing the SVM ourselves, we will gain a better understanding of its inner workings. This hands-on experience will help clarify any confusion and solidify our knowledge of SVMs. If you have any questions, we recommend holding them until the next material, where we will provide a recap and simplify the mathematical concepts we have covered so far.

Feel free to leave your comments and questions below. Thank you for watching and stay tuned for the next material where we will delve deeper into Support Vector Machines.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SUPPORT VECTOR MACHINE OPTIMIZATION**

The Support Vector Machine (SVM) is a powerful classification technique in machine learning. To fully comprehend its implementation and optimization, it is essential to understand the fundamental concepts and mathematical formulations that underpin it.

The equation for a hyperplane, which is central to SVM, is given by:

$$\mathbf{x} \cdot \mathbf{w} + b = 0$$

Here,  $\mathbf{x}$  represents the feature vector,  $\mathbf{w}$  is the weight vector, and  $b$  is the bias term.

For a positive class support vector, the equation is:

$$\mathbf{x}_i \cdot \mathbf{w} + b = 1$$

For a negative class support vector, it is:

$$\mathbf{x}_i \cdot \mathbf{w} + b = -1$$

These equations are specific to the support vectors, which are the data points that lie closest to the decision boundary. The decision boundary itself is defined by:

$$\mathbf{x} \cdot \mathbf{w} + b = 0$$

When considering a point  $\mathbf{x}_i$  such that:

$$\mathbf{x}_i \cdot \mathbf{w} + b = 0.98$$

This point lies between the support vector hyperplane and the decision boundary. If the value is positive but less than 1, it indicates that the point is within the margin but still classified as part of the positive class.

The classification of a feature set after training the SVM classifier is determined by the sign of the decision function:

$$\text{sign}(\mathbf{x}_i \cdot \mathbf{w} + b)$$

If the result is positive, the point belongs to the positive class; if negative, it belongs to the negative class. If the result is zero, the point lies on the decision boundary.

The optimization of the SVM involves finding the optimal  $\mathbf{w}$  and  $b$ . The objective is to minimize the magnitude of the weight vector  $\mathbf{w}$ , which can be expressed as:

$$\|\mathbf{w}\| = \sqrt{\sum_j w_j^2}$$

This is analogous to the Euclidean norm or the length of the vector in multi-dimensional space.

The optimization is subject to the constraint:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1$$

for all training samples  $i$ , where  $y_i$  is the class label (+1 or -1).

To summarize, the SVM optimization problem aims to find the weight vector  $\mathbf{w}$  and bias  $b$  that minimize the norm of  $\mathbf{w}$  while satisfying the classification constraints. This can be formulated as a constrained optimization problem:

Minimize:

$$\frac{1}{2} \|\mathbf{w}\|^2$$

Subject to:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1$$

This problem can be solved using various optimization techniques, such as quadratic programming. Once the optimal  $\mathbf{w}$  and  $b$  are determined, the classification of new data points becomes straightforward using the decision function.

In practical implementation using Python, libraries such as scikit-learn provide built-in functions to perform SVM classification, making it easier to apply SVMs to real-world datasets.

In the context of support vector machines (SVMs), the classification problem can be expressed mathematically. The class variable, denoted as  $y$ , takes values of either  $+1$  or  $-1$ . The objective is to find a hyperplane that best separates the data points of different classes. This can be formulated as:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$

where  $\mathbf{w}$  is the weight vector,  $\mathbf{x}_i$  represents the feature vector of the  $i$ -th sample, and  $b$  is the bias term. The goal is to find the vector  $\mathbf{w}$  with the smallest magnitude, which translates to minimizing  $\|\mathbf{w}\|$ , while ensuring that the constraint  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$  is satisfied for all training samples.

To solve this optimization problem, one must minimize the magnitude of  $\mathbf{w}$ , or more formally, minimize  $\frac{1}{2} \|\mathbf{w}\|^2$ , subject to the constraints. This is a quadratic programming problem, characterized by a convex optimization landscape. Convex problems are advantageous because they have a single global minimum, making them easier to solve compared to non-convex problems.



Visualizing the convex optimization landscape, one can think of it as a bowl-shaped curve. If a ball is placed on the surface of this bowl, it will roll down to the lowest point, representing the minimum value of the objective function. In mathematical terms, this means that the optimization algorithm will converge to the global minimum of  $\|\mathbf{w}\|$ .

Given a set of known features  $\mathbf{x}_i$ , the task is to find the optimal  $\mathbf{w}$  and  $b$ . Consider an example where the feature vectors are two-dimensional, i.e.,  $\mathbf{w}$  is a 1x2 matrix. An example of  $\mathbf{w}$  could be  $[5, 3]$ . The norm (or magnitude) of  $\mathbf{w}$  is calculated as:

$$\|\mathbf{w}\| = \sqrt{5^2 + 3^2} = \sqrt{34}$$

If  $\mathbf{w}$  were  $[-5, 3]$ , the norm remains  $\sqrt{34}$ . However, the sign of the components of  $\mathbf{w}$  affects the dot product  $\mathbf{x}_i \cdot \mathbf{w}$ . For instance, if  $\mathbf{x}_i = [1, 2]$  and  $\mathbf{w} = [5, 3]$ , the dot product would be:

$$\mathbf{x}_i \cdot \mathbf{w} = 1 \cdot 5 + 2 \cdot 3 = 11$$

If  $\mathbf{w} = [-5, 3]$ , the dot product changes to:

$$\mathbf{x}_i \cdot \mathbf{w} = 1 \cdot (-5) + 2 \cdot 3 = 1$$

This demonstrates that the sign of the components of  $\mathbf{w}$  significantly impacts the outcome of the dot product, and thus the classification result.

The optimization process in SVMs involves finding the weight vector  $\mathbf{w}$  with the smallest magnitude that satisfies the classification constraints. This is achieved by solving a convex quadratic programming problem, ensuring a unique and optimal solution.

Support Vector Machine (SVM) optimization involves finding the optimal hyperplane that separates different classes in a dataset. Traditionally, we have been dealing with problems where the features allow for a straightforward separation with a hyperplane of a certain slope. However, when the feature set changes, the direction and orientation of the hyperplane must also adapt.

For instance, consider a scenario where the initial vector  $\mathbf{w}$  is set to  $[5, 5]$ . The goal is to find a bias  $b$  such that the inequality  $y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1$  holds true for all training samples  $(\mathbf{x}_i, y_i)$ . If such a  $b$  exists, we record the magnitude and direction of  $\mathbf{w}$ . We then decrement  $\mathbf{w}$  to  $[4, 4]$  and repeat the process. This iterative stepping down of  $\mathbf{w}$  continues until no further decrement is possible without violating the inequality.

During each iteration, it is essential to test not only the current vector  $\mathbf{w}$ , but also its negative and permutations, such as  $[-5, 5]$ ,  $[5, -5]$ , and  $[-5, -5]$ . This comprehensive testing ensures that no potential optimal vector is overlooked.

The dot product plays a important role in these calculations. For instance, if  $\mathbf{w} = [5, 5]$ , we compute the dot product with vectors like  $[-1, 1]$ ,  $[1, -1]$ , and  $[-1, -1]$ . Each step involves evaluating four different configurations, making the process computationally intensive.

The convex nature of the optimization problem simplifies the search for the minimum value. By taking incremental steps, one can avoid recalculating values for points that are already known to be suboptimal. If

larger steps are taken and the minimum value is overshoot, backtracking is necessary to refine the search. This iterative process of stepping forward and back ensures convergence to the optimal solution.

The optimization of SVM involves iterative adjustments of the weight vector  $\mathbf{w}$  and comprehensive testing of its permutations to ensure the best separation of classes. The convexity of the problem aids in efficiently finding the minimum value, making SVM a robust method for classification tasks.

In the context of support vector machine (SVM) optimization, the goal is to find the optimal separating hyperplane that maximizes the margin between different classes in the feature space. The optimization process involves finding the global minimum of a convex problem, ensuring that the solution is not trapped in local minima.

To begin, consider initializing the weight vector  $\mathbf{w}$  at a starting value, say  $\mathbf{w} = (10, 10)$ . This vector will be iteratively adjusted to satisfy the SVM constraints. The primary constraint for SVM is given by:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1$$

where  $y_i$  is the class label,  $\mathbf{x}_i$  is the feature vector,  $\mathbf{w}$  is the weight vector, and  $b$  is the bias term. The objective is to find  $\mathbf{w}$  and  $b$  such that this inequality holds for all training samples.

The optimization process involves iteratively adjusting  $\mathbf{w}$  and  $b$ . Initially,  $\mathbf{w}$  is plugged into the equation, and  $b$  is incrementally adjusted to find a value that satisfies the constraint. Each combination of  $\mathbf{w}$  and  $b$  that meets the condition is stored in a dictionary, where the key is the magnitude of  $\mathbf{w}$  and the value is the tuple  $(\mathbf{w}, b)$ .

The magnitude of  $\mathbf{w}$  is given by:

$$\text{magnitude}(\mathbf{w}) = \|\mathbf{w}\| = \sqrt{w_1^2 + w_2^2 + \dots + w_n^2}$$

The dictionary is used to keep track of all valid  $\mathbf{w}$  and  $b$  pairs, and the pair with the smallest magnitude is selected as the optimal solution. This process ensures that the solution corresponds to the maximum margin hyperplane.

The iterative optimization process can be visualized as taking steps towards the minimum point of a convex function. If the steps are too large, the solution might overshoot the minimum, necessitating smaller steps to fine-tune the solution. This process continues until the global minimum is reached or approximated closely.

In SVM, the problem is convex, meaning there is no risk of getting trapped in local minima. This ensures that the global minimum found is indeed the optimal solution.

To verify the correctness of the solution, one can plot the decision boundary, support vectors, and margins. The support vectors are the data points that lie closest to the decision boundary and satisfy the equation:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) = 1$$

During the optimization, if the support vectors' values are very close to 1 (e.g., 1.001), it indicates that the solution is near optimal. The presence of at least one data point in each class with a value close to 1 confirms the identification of support vectors and the correctness of the optimization.

The SVM optimization process involves initializing  $\mathbf{w}$ , iteratively adjusting  $\mathbf{w}$  and  $b$  to satisfy the SVM constraints, storing valid solutions, and selecting the one with the smallest magnitude. The convex nature of the problem ensures that the global minimum found is the optimal solution, which can be verified by the proximity of support vector values to 1.

To effectively determine whether one has closely approximated an optimal solution in Support Vector Machine (SVM) optimization, it is important to observe the outputs from the algorithm. When the outputs for both positive and negative feature sets are very close to one, it indicates proximity to an optimal solution. This is a significant marker in the iterative process of optimization, guiding whether to continue taking smaller steps or to accept the minimum found.

Optimization is a vast field, encompassing much more than machine learning. Specifically, SVM optimization is a subset of convex optimization problems. Convex optimization involves finding a minimum of a convex function over a convex set, ensuring that any local minimum is also a global minimum.

Python provides several modules to assist with convex optimization. One notable library is CVXOpt, which is designed to solve quadratic programming (QP) problems and includes specific algorithms for SVM. Another useful module is LibSVM, which is also employed in the Scikit-learn library for SVM optimization tasks. However, for educational purposes, it is beneficial to manually implement the optimization process to gain a deeper understanding of the underlying mechanics.

To illustrate, consider the following simplified example of SVM optimization:

1. **Objective Function:** The goal is to minimize the function:

$$\frac{1}{2} \|\mathbf{w}\|^2$$

subject to the constraints:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad \text{for all } i.$$

2. **Lagrange Multipliers:** Introduce Lagrange multipliers  $\alpha_i$  for each constraint:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i \alpha_i [y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1].$$

3. **Optimization:** Solve the dual problem:

$$\max_{\alpha} \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (\mathbf{x}_i \cdot \mathbf{x}_j)$$

subject to:

$$\sum_i \alpha_i y_i = 0 \quad \text{and} \quad \alpha_i \geq 0.$$

4. **Solution:** The optimal weights  $\mathbf{w}$  can be derived from:

$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i.$$

5. **Decision Boundary:** The decision boundary is given by:

$$\mathbf{w} \cdot \mathbf{x} + b = 0.$$

By manually coding these steps, one gains an intuitive understanding of the SVM optimization process. This hands-on approach demystifies the complex mathematical formulations and enhances comprehension through practical implementation.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: CREATING AN SVM FROM SCRATCH**

In this didactic material, we will be discussing the support vector machine (SVM) algorithm in the context of machine learning. The SVM is a powerful supervised learning algorithm used for classification and regression tasks. In this tutorial, we will learn how to create an SVM from scratch using Python.

Before we begin, it is important to note that we have already covered the theory and logic behind the SVM in previous parts of this tutorial series. If you are unfamiliar with the SVM or need a refresher, we recommend going back and reviewing the relevant material.

To get started, we need to import the necessary libraries. We will be using Matplotlib for data visualization, so we import it as plt. Additionally, we import the NumPy library as np for numerical operations. These libraries are essential for creating and visualizing our SVM.

Next, we will create some simple basic data for our SVM. We define a dictionary called "data\_dict" with two keys: -1 and 1. Each key corresponds to a class, and the values are lists of lists representing data points. We populate the lists with some sample data points.

Once we have our data, we can proceed to build our support vector machine class. In object-oriented programming, classes are used to define objects with their own properties and methods. Our SVM class will allow us to train the model, make predictions, and visualize the results.

We define the class as "support vector machine" and create an initialization method, denoted as "init". The init method is automatically called when we create an instance of the class. Within the init method, we set the "visualization" attribute to True by default. This attribute determines whether we want to visualize the data and results. We also define the "colors" attribute, which specifies the colors for different classes when visualizing the data.

If the "visualization" attribute is set to True, we create a figure using plt.figure() and assign it to the "fig" attribute. We also create axes using self.fig.add\_subplot() and assign them to the "axes" attribute. These steps are necessary for plotting the data.

Now that we have set up the basic structure of our SVM class, we can move on to implementing the SVM algorithm itself. However, since the provided transcript is incomplete, we will skip this part.

To summarize, in this tutorial, we have discussed the support vector machine algorithm and how to create an SVM from scratch using Python. We have covered the necessary libraries, data preparation, and the initialization of our SVM class. In the next parts of this tutorial series, we will continue building the SVM and explore its functionality.

In this didactic material, we will discuss the process of creating a Support Vector Machine (SVM) from scratch using Python. SVM is a powerful machine learning algorithm used for classification tasks. We will cover the steps involved in creating an SVM and explain the concepts along the way.

To begin, we need to import the necessary libraries. One of the libraries we will be using is Matplotlib, which is a popular data visualization library in Python. Matplotlib allows us to create plots and graphs to visualize our data.

Next, we will create a subplot in Matplotlib. A subplot is a grid of plots within a single figure. In this case, we are creating a subplot with a one-by-one grid and plot number one. This subplot will be used for visualization purposes later on.

Moving on to the main part of our SVM, we need to define the 'fit' method. The 'fit' method is responsible for training our SVM model. It takes in the 'self' parameter, which allows us to share variables and data within the class. Additionally, it takes in the 'data' parameter, which represents the input data for training.

Inside the 'fit' method, we will perform the necessary calculations to optimize our SVM. However, in this initial

stage, we will simply use the 'pass' statement to indicate that no further actions are required.

Another important method in our SVM is the 'predict' method. The 'predict' method is used to make predictions on new, unseen data. Similar to the 'fit' method, it takes in the 'self' parameter and the 'data' parameter, representing the input data for prediction.

Inside the 'predict' method, we will calculate the classification for each data point. To do this, we will use the formula:  $\text{classification} = \text{sign}(X \text{ dot } W + B)$ , where 'X' represents the input data, 'W' represents the weight vector, 'B' represents the bias term, and 'sign' is a function that returns the sign of a number.

To implement this formula in Python, we will use the NumPy library. NumPy provides a 'sign' function that allows us to calculate the sign of a number. We will use the 'dot' function in NumPy to perform the dot product between 'X' and 'W'.

Once we have the classification for each data point, we will return the classification as the output of the 'predict' method.

At this point, we have covered the basic structure of our SVM. However, there are still some missing components, such as the weight vector 'W' and the bias term 'B'. These values will be optimized and set during the 'fit' method, which we will cover in a future tutorial.

In the next tutorial, we will consider the optimization process and discuss how to find the optimal values for 'W' and 'B'. We will also cover the visualization of our SVM using Matplotlib.

If you have any questions or concerns up to this point, please feel free to reach out. Thank you for watching and for your continued support and subscriptions.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SVM TRAINING**

Support Vector Machines (SVM) is a popular machine learning algorithm used for classification and regression tasks. In this tutorial, we will focus on SVM training and optimization.

The SVM algorithm aims to find the best hyperplane that separates different classes in the data. The hyperplane is defined by a weight vector ( $W$ ) and a bias term ( $B$ ). The goal is to find the optimal values for  $W$  and  $B$  that maximize the margin between the classes.

To train an SVM model, we need to solve an optimization problem. The objective is to minimize the loss function while satisfying certain constraints. In SVM, the loss function is based on hinge loss, which penalizes misclassified samples. The optimization problem is a convex problem, meaning it has a unique global minimum.

The training process involves finding the optimal values for  $W$  and  $B$  that minimize the loss function. This is done using an optimization technique called Sequential Minimal Optimization (SMO). SMO is a rudimentary method for solving the optimization problem, and it becomes computationally expensive as the dataset size increases.

To better understand optimization, you can refer to the Stanford Convex Optimization book, which provides in-depth knowledge on the subject. Additionally, the CVX Optima Joule website offers resources on convex optimization, including research papers and source code.

In SVM training, we start by loading the dataset. The dataset is stored in the "data" variable. We then define an empty dictionary called "opt\_dict" to store the magnitude of  $W$  as the key and the corresponding values of  $W$  and  $B$  as the list. This dictionary will be populated during the training process.

Next, we define the transforms. These transforms are applied to the weight vector ( $W$ ) at each step of the training process. The transforms help in handling both positive and negative values of  $W$ . By applying these transforms, we ensure that the dot product of  $W$  and the feature set is correctly calculated.

Once the dataset is loaded and the transforms are defined, we can proceed with the optimization process. The goal is to find the optimal values for  $W$  and  $B$  that minimize the loss function. This is achieved by iteratively updating the values of  $W$  and  $B$  using the SMO algorithm.

It's important to note that the training process becomes computationally expensive as the dataset size increases. Each data sample needs to be checked to determine if it satisfies the classification condition. This process is repeated for every data sample, making it time-consuming for large datasets.

SVM training involves solving an optimization problem to find the best hyperplane that separates different classes in the data. The training process uses the SMO algorithm and can be computationally expensive for large datasets. Understanding optimization techniques and convex optimization can further enhance your knowledge in this area.

In support vector machine (SVM) training, it is important to determine the maximum and minimum ranges for our graph and the initial values for the variables  $W$  and  $B$ . To achieve this, we can write a function or a for loop to iterate through the data. We iterate through the classes and feature sets, and for each feature, we append it to a list called "all\_data". We then find the maximum and minimum values in this list to obtain the maximum and minimum feature values. Finally, we set "all\_data" to None to free up memory.

To determine the step sizes, we consider the concept of taking big steps initially and gradually reducing the step size. We start with a step size of 0.1 times the maximum feature value. This corresponds to the big steps. Once we find a good value, we reduce the step size to 1% of the maximum feature value. And after finding an even better value, we take even smaller steps. However, it is important to note that taking smaller steps beyond a certain point becomes expensive.

To improve the SVM optimization problem, there are a few areas to consider. Firstly, we can explore the

possibility of threading or multiprocessing to run the step functions simultaneously. This can potentially speed up the process. Additionally, we can tweak various parameters to achieve a balance between accuracy and efficiency.

Another variable to set is the "V range multiple" which is set to 5. This variable determines the step size for B. It is important to note that B does not need to be as precise as W, and making it more precise can be expensive.

In SVM training, we determine the maximum and minimum feature values to establish the ranges for our graph and initial variable values. We then take steps of varying sizes, starting with big steps and gradually reducing the step size. We can explore threading or multiprocessing to improve efficiency. Additionally, we set the "V range multiple" to control the step size for B. It is important to find a balance between accuracy and computational cost.

Support Vector Machines (SVM) is a popular machine learning algorithm used for classification and regression tasks. In this didactic material, we will focus on SVM training using Python.

Before diving into the details, let's briefly discuss the concept of SVM. SVM is a supervised learning algorithm that analyzes data and builds a classification model. It is based on the idea of finding a hyperplane that separates the data into different classes, maximizing the margin between them.

In SVM training, we start by initializing some parameters. For example, we set the value of 'B' to 5 and 'latest\_optimum' to 'self.Max\_feature\_value \* 10'. The 'latest\_optimum' represents the starting basic list, where the first element in vector W is multiplied by the same value for every element in the vector.

Once the initial values are set, we proceed with the stepping process. We iterate over a range of step sizes and perform the following steps for each step:

1. We initialize vector W using the latest optimum value.
2. We set the 'optimized' value to False.
3. We check if the optimization is complete by evaluating the 'optimized' value. If it is False, we continue with the optimization process.

The optimization process continues until there are no more steps to take. This is possible because SVM is a convex problem, meaning there is a single global minimum.

In the next tutorial, we will explore the remaining steps to complete the SVM training process. We will also discuss how to graph the results and finalize the SVM implementation.

If you have any questions or concerns, please feel free to leave them in the comments section. Thank you for watching and for your support!



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SVM OPTIMIZATION**

In this tutorial, we will continue building our support vector machine (SVM) class and focus on the optimization process. Specifically, we will discuss the iteration through B values to find the maximum bias possible.

To start, we will iterate through the B values using the Numpy function `np.arange()`. This function allows us to specify the step size for iteration. In this case, we want the step size to be the product of the maximum feature value and a range multiple. The range multiple is set to -1 times the maximum feature value, resulting in a negative range that starts from a value like -50.

Next, we need to define the step size for each iteration. Since the optimization for B is more computationally expensive than that for W, we can take larger steps. Therefore, we multiply the step size by a factor of 5.

Once we have set up the iteration for B, we can proceed to transform the W values. For each transformation, we apply the transformation function to the original W value. The transformation function is defined as `W_transformed = W * transformation`.

Now that we have the transformed W values, we can move on to testing. We start by assuming that we have found the optimal option and set `found_option` to True. Then, we iterate through the data to check if the SVM fits the data correctly. This is done by iterating through each class and each data point within that class.

Finally, we can conclude that the weakest link in the SVM is the need to run calculations on all the data to ensure it fits correctly. This can be computationally expensive, especially when dealing with large datasets. However, there are methods such as Sequential Minimal Optimization (SMO) that attempt to mitigate this issue to some extent.

In this tutorial, we have discussed the optimization process for the support vector machine, specifically focusing on the iteration through B values and the testing of the SVM on the data.

Support Vector Machine (SVM) optimization is an important aspect of machine learning with Python. In SVM, the optimization process involves finding the best hyperplane that separates the data points into different classes. This didactic material will explain the steps involved in SVM optimization using Python.

The first step in SVM optimization is defining the constraint function. The constraint function is given by the equation  $y_i \cdot (w \cdot x_i + b) \geq 1$ . This equation ensures that the data points are correctly classified. If this constraint is not satisfied for any sample in the dataset, the optimization process is considered unsuccessful.

To implement SVM optimization in Python, we start by initializing a variable called `found_option` as true. We iterate through the dataset and check if the constraint function is satisfied for each sample. If we find a sample that does not satisfy the constraint, we set `found_option` to false and break the loop. This optimization technique allows us to stop checking further samples from that class or the other class, thus improving efficiency.

Next, we perform the transformation separately for each option. If the constraint function is satisfied for all samples, we proceed with the optimization process. We calculate the magnitude of the vector using the equation  $\sqrt{a^2 + b^2 + c^2 + \dots}$ . This magnitude is then assigned to the variable `opt_dict` using the W and B values.

After running through all the options, we check if the W value is less than zero. If it is, we set the variable `optimized` as true, indicating that the optimization step is successful. Otherwise, we take the step by subtracting the `step` value from the W vector.

Finally, we sort the magnitudes of the `opt_dict` values in ascending order and assign the first element to the variable `opt_choice`. We then set the values of W and B as `opt_choice[0]` and `opt_choice[1]` respectively.

By following these steps, we can optimize the Support Vector Machine model using Python. This optimization process helps in finding the best hyperplane that maximally separates the data points of different classes.

In the previous material, we discussed the optimization process in Support Vector Machines (SVM) for machine learning using Python. We explored the concept of norms and how they relate to finding the optimal choice in SVM.

To recap, norms represent the magnitudes of vectors. In the optimization process, we sort the norms in ascending order to find the smallest norm, which will be our optimal choice. We store this optimal choice in a dictionary, where the key is the magnitude of the vector and the value is the vector itself.

Additionally, we set the latest optimum value by adding the step size to the zeroeth element of the optimal choice. This allows us to reference the latest optimum value as we go through each step of the optimization process. By modifying this value with each step, we can make our predictions more precise.

To determine whether we need to take another step in the optimization process, we use the support vectors. Support vectors are calculated by taking the dot product of the vector with the weight vector ( $W$ ) and adding the bias ( $B$ ). If the value is close to 1 in both the positive and negative classes, we have found a good value for  $W$  and  $B$ . The closeness to 1 can be defined based on the problem at hand. If the value is not close to 1, we continue taking steps, with each step size being an order of magnitude smaller than the previous one.

It's important to note that not all problems can be optimized using a linear kernel. In some cases, it may be necessary to stop the optimization process if a satisfactory solution cannot be found.

At this point, we have completed the optimization algorithm for SVM. In the next tutorial, we will consider predicting outcomes using our optimized model. We will also cover the visualization of our results.

If you have any questions or concerns regarding the material covered so far, please feel free to leave them in the comments section. Thank you for your support and subscriptions, and stay tuned for the next tutorial!

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: COMPLETING SVM FROM SCRATCH**

To implement a Support Vector Machine (SVM) from scratch in Python, one must understand both the theoretical foundation and practical coding techniques. This material will guide you through completing an SVM, including optimization, prediction, and visualization.

First, consider the prediction method. This method assesses new data points to determine their classification based on the trained SVM model. If visualization is enabled, the data points are plotted to visually verify the classification.

1.	<code>def predict(self, features):</code>
2.	<code>    classification = np.sign(np.dot(np.array(features), self.w) + self.b)</code>
3.	<code>    if classification != 0 and self.visualization:</code>
4.	<code>        self.ax.scatter(features[0], features[1], s=200, marker='*', c=self.colors[c</code>
	<code>lassification])</code>
5.	<code>    return classification</code>

In the `predict` method, the classification is determined by the sign of the dot product of the input features and the weight vector `w`, added to the bias `b`. If visualization is enabled, the method plots the data point using `matplotlib`.

Next, the `visualize` method is designed to plot the data points and the hyperplanes that define the SVM decision boundary. This is important for understanding how the SVM separates the data.

1.	<code>def visualize(self):</code>
2.	<code>    [[self.ax.scatter(x[0], x[1], s=100, color=self.colors[i]) for x in data_dict[i]</code>
	<code>] for i in data_dict]</code>
3.	
4.	<code>    def hyperplane(x, w, b, v):</code>
5.	<code>        return (-w[0] * x - b + v) / w[1]</code>
6.	
7.	<code>    datarange = (self.min_feature_value * 0.9, self.max_feature_value * 1.1)</code>
8.	<code>    hyp_x_min = datarange[0]</code>
9.	<code>    hyp_x_max = datarange[1]</code>
10.	
11.	<code>    # Positive support vector hyperplane</code>
12.	<code>    psv1 = hyperplane(hyp_x_min, self.w, self.b, 1)</code>
13.	<code>    psv2 = hyperplane(hyp_x_max, self.w, self.b, 1)</code>
14.	<code>    self.ax.plot([hyp_x_min, hyp_x_max], [psv1, psv2], 'k')</code>
15.	
16.	<code>    # Negative support vector hyperplane</code>
17.	<code>    nsv1 = hyperplane(hyp_x_min, self.w, self.b, -1)</code>
18.	<code>    nsv2 = hyperplane(hyp_x_max, self.w, self.b, -1)</code>
19.	<code>    self.ax.plot([hyp_x_min, hyp_x_max], [nsv1, nsv2], 'k')</code>
20.	
21.	<code>    # Decision boundary hyperplane</code>
22.	<code>    db1 = hyperplane(hyp_x_min, self.w, self.b, 0)</code>
23.	<code>    db2 = hyperplane(hyp_x_max, self.w, self.b, 0)</code>
24.	<code>    self.ax.plot([hyp_x_min, hyp_x_max], [db1, db2], 'y--')</code>
25.	
26.	<code>    plt.show()</code>

In the `visualize` method, the data points are first plotted. The `hyperplane` function is defined to calculate the hyperplane values. The hyperplane is defined by the equation:

$$v = x \cdot w + b$$

To find specific hyperplanes (support vectors and decision boundary), the function rearranges this equation to solve for  $x_2$ :

$$x_2 = \frac{-w_0 \cdot x_1 - b + v}{w_1}$$

Where  $v$  represents the value for which we are solving (1 for positive support vector, -1 for negative support vector, and 0 for the decision boundary).

The visualization includes plotting:

1. Positive support vector hyperplane.
2. Negative support vector hyperplane.
3. Decision boundary hyperplane.

These hyperplanes help in understanding how the SVM classifies data points.

Completing an SVM from scratch involves defining methods for prediction and visualization, which rely on the mathematical foundation of hyperplanes and dot products. The implementation in Python leverages libraries such as NumPy for mathematical operations and Matplotlib for plotting.

In the context of machine learning, a Support Vector Machine (SVM) is a powerful supervised learning algorithm used for classification and regression tasks. The primary objective of an SVM is to find the optimal hyperplane that best separates the data points of different classes. The mathematical foundation of an SVM involves finding the weights  $\mathbf{w}$  and bias  $b$  that define this hyperplane.

Given a set of training data points  $(\mathbf{x}_i, y_i)$ , where  $\mathbf{x}_i$  represents the feature vectors and  $y_i$  the class labels, the SVM aims to solve the optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to the constraints:

$$y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 \quad \forall i$$

This ensures that the data points are correctly classified with a margin of at least 1. The decision boundary is defined by the equation:

$$\mathbf{w} \cdot \mathbf{x} + b = 0$$

To visualize the SVM, one can plot the decision boundary along with the support vectors. However, it is important to note that this visualization is for human interpretation and has no bearing on the actual SVM algorithm. The visualization typically involves plotting the decision boundary and the margins, which are hyperplanes defined by:

$$\mathbf{w} \cdot \mathbf{x} + b = \pm 1$$

To implement this in Python, we can use the following steps:

### 1. Define the Data Range:

To ensure the graph has some extra space around the data points, we define the data range with a margin of 10%.

```
1. data_range = (self.min_feature_value * 0.9, self.max_feature_value * 1.1)
```

### 2. Create the Hyperplanes:

The hyperplanes are defined by the equations  $\mathbf{w} \cdot \mathbf{x} + b = 1$  (positive support vector),  $\mathbf{w} \cdot \mathbf{x} + b = -1$  (negative support vector), and  $\mathbf{w} \cdot \mathbf{x} + b = 0$  (decision boundary).

```
1. def hyperplane(x, w, b, v):
2.     return (-w[0] * x - b + v) / w[1]
3.
4. hyperplane_x_min = data_range[0]
5. hyperplane_x_max = data_range[1]
6.
7. # Positive support vector hyperplane
8. psv1 = hyperplane(hyperplane_x_min, self.w, self.b, 1)
9. psv2 = hyperplane(hyperplane_x_max, self.w, self.b, 1)
10.
11. # Negative support vector hyperplane
12. nsv1 = hyperplane(hyperplane_x_min, self.w, self.b, -1)
13. nsv2 = hyperplane(hyperplane_x_max, self.w, self.b, -1)
14.
15. # Decision boundary hyperplane
16. db1 = hyperplane(hyperplane_x_min, self.w, self.b, 0)
17. db2 = hyperplane(hyperplane_x_max, self.w, self.b, 0)
```

### 3. Plot the Hyperplanes:

Using a plotting library such as Matplotlib, the hyperplanes can be visualized.

```
1. import matplotlib.pyplot as plt
2.
3. plt.plot([hyperplane_x_min, hyperplane_x_max], [psv1, psv2], 'k')
4. plt.plot([hyperplane_x_min, hyperplane_x_max], [nsv1, nsv2], 'k')
5. plt.plot([hyperplane_x_min, hyperplane_x_max], [db1, db2], 'g--')
6.
7. plt.scatter(data[:, 0], data[:, 1], c=labels)
8. plt.show()
```

While the visualization of the decision boundary and support vectors is useful for understanding and interpreting the SVM, the core of the algorithm focuses on finding the optimal weights  $\mathbf{w}$  and bias  $b$  that maximize the margin between the classes. The hyperplanes defined by  $\mathbf{w} \cdot \mathbf{x} + b = \pm 1$  are critical in determining the support vectors, which are the data points closest to the decision boundary.

In the context of machine learning, a Support Vector Machine (SVM) is a powerful supervised learning model used for classification and regression tasks. The primary objective of an SVM is to find the optimal hyperplane that best separates the data into different classes. This hyperplane is defined by support vectors, which are the data points closest to the decision boundary.

To implement an SVM from scratch in Python, one must follow several key steps: initialization, optimization (fitment), prediction, and visualization. Below is a structured approach to building and visualizing an SVM.

## INITIALIZATION

First, define the SVM class. This class will include methods for fitting the model to the data, making predictions, and visualizing the results.

1.	class SupportVectorMachine:
2.	def __init__(self):
3.	self.w = None
4.	self.b = None
5.	
6.	def fit(self, data):
7.	# Implementation of the optimization algorithm to find w and b
8.	pass
9.	
10.	def predict(self, features):
11.	# Return the prediction based on the sign of (XW + b)
12.	return np.sign(np.dot(np.array(features), self.w) + self.b)
13.	
14.	def visualize(self):
15.	# Visualization code to plot the decision boundary and support vectors
16.	pass

### OPTIMIZATION (FITMENT)

The optimization process involves finding the weight vector `w` and bias `b` that maximize the margin between the different classes. This is typically done using quadratic programming or gradient descent methods.

### PREDICTION

The prediction function uses the sign of the decision function  $f(x) = x \cdot w + b$  to classify new data points.

### VISUALIZATION

Visualization helps in understanding how well the SVM has separated the classes. The decision boundary, support vectors, and margins can be plotted for better insight.

1.	import numpy as np
2.	import matplotlib.pyplot as plt
3.	
4.	class SupportVectorMachine:
5.	def __init__(self):
6.	self.w = None
7.	self.b = None
8.	
9.	def fit(self, data):
10.	# Simplified example of fitting the model
11.	self.w = np.array([1, 1])
12.	self.b = -1
13.	
14.	def predict(self, features):
15.	return np.sign(np.dot(np.array(features), self.w) + self.b)
16.	
17.	def visualize(self, data):
18.	def hyperplane(x, w, b, offset):
19.	return (-w[0] * x - b + offset) / w[1]
20.	
21.	plt.scatter(data['pos'][:, 0], data['pos'][:, 1], color='b')
22.	plt.scatter(data['neg'][:, 0], data['neg'][:, 1], color='r')
23.	
24.	x_min, x_max = -10, 10
25.	plt.plot([x_min, x_max], [hyperplane(x_min, self.w, self.b, 0), hyperplane(x_max, self.w, self.b, 0)], 'k')
26.	plt.plot([x_min, x_max], [hyperplane(x_min, self.w, self.b, 1), hyperplane(x_max, self.w, self.b, 1)], 'k--')
27.	plt.plot([x_min, x_max], [hyperplane(x_min, self.w, self.b, -1), hyperplane(x_max, self.w, self.b, -1)], 'k--')
28.	

29.	<code>plt.show()</code>
30.	
31.	<code># Example usage</code>
32.	<code>data = {</code>
33.	<code>    'pos': np.array([[3, 4], [1, 2], [3, 3]]),</code>
34.	<code>    'neg': np.array([[6, 7], [8, 9], [7, 8]])</code>
35.	<code>}</code>
36.	
37.	<code>svm = SupportVectorMachine()</code>
38.	<code>svm.fit(data)</code>
39.	<code>svm.visualize(data)</code>

In the example above:

1. `data` is a dictionary containing positive and negative samples.
2. The `fit` method is simplified for illustrative purposes.
3. The `visualize` method plots the decision boundary and margins.

The decision boundary is represented by the hyperplane where the decision function equals zero. The margins are plotted as dashed lines, indicating the boundaries within which the support vectors lie.

### MATHEMATICAL FORMULATION

The SVM optimization problem can be formulated as follows:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$$

subject to:

$$y_i(\mathbf{x}_i \cdot \mathbf{w} + b) \geq 1, \quad \forall i$$

Where:

- $\mathbf{w}$  is the weight vector.
- $b$  is the bias.
- $\mathbf{x}_i$  are the input features.
- $y_i$  are the class labels (either 1 or -1).

This formulation seeks to maximize the margin while ensuring that all data points are correctly classified. Building an SVM from scratch involves understanding the mathematical foundations, implementing the optimization algorithm, and effectively visualizing the results. This approach provides deep insights into the inner workings of SVMs and enhances one's ability to customize and extend the model for specific applications.

In the context of machine learning, Support Vector Machines (SVM) are a powerful tool for classification tasks. This material will guide you through the process of implementing SVM from scratch using Python, focusing on the prediction phase and visualization of results.

To begin with, after training the SVM model, the next step involves making predictions on new data points. Suppose we have a list of new data points for which we want to predict the class. We initialize this list as follows:

1.	<code>predict_us = [</code>
2.	<code>    [0, 1],</code>
3.	<code>    [1, 3],</code>
4.	<code>    [3, 4],</code>
5.	<code>    [3, 5],</code>
6.	<code>    [5, 5],</code>

## EUROPEAN IT CERTIFICATION CURRICULUM SELF-LEARNING PREPARATORY MATERIALS

7.	[5, 6],
8.	[6, -5],
9.	[5, 8]
10.	]

We then use the `predict` method of our SVM class to classify these points. The prediction process involves calculating the dot product of the input data point with the weight vector `w` and adding the bias term `b`. The sign of the result determines the class:

1.	for p in predict_us:
2.	print(svm.predict(p))

The visualization of the results can be achieved by plotting the data points along with the decision boundary. The decision boundary is defined by the equation  $w \cdot x + b = 0$ . Points on one side of the boundary belong to one class, while points on the other side belong to the opposite class.

The following code snippet demonstrates how to visualize the data points and the decision boundary:

1.	import matplotlib.pyplot as plt
2.	
3.	# Plotting the support vectors
4.	for p in predict_us:
5.	plt.scatter(p[0], p[1], s=100, marker='*')
6.	
7.	# Assuming svm is an instance of the SVM class
8.	# Plotting the decision boundary
9.	def plot_decision_boundary(svm):
10.	def hyperplane(x, w, b, v):
11.	return (-w[0] * x - b + v) / w[1]
12.	
13.	datarange = (0, 10)
14.	hyp_x_min = datarange[0]
15.	hyp_x_max = datarange[1]
16.	
17.	# Positive support vector hyperplane
18.	psv1 = hyperplane(hyp_x_min, svm.w, svm.b, 1)
19.	psv2 = hyperplane(hyp_x_max, svm.w, svm.b, 1)
20.	plt.plot([hyp_x_min, hyp_x_max], [psv1, psv2], 'k')
21.	
22.	# Negative support vector hyperplane
23.	nsv1 = hyperplane(hyp_x_min, svm.w, svm.b, -1)
24.	nsv2 = hyperplane(hyp_x_max, svm.w, svm.b, -1)
25.	plt.plot([hyp_x_min, hyp_x_max], [nsv1, nsv2], 'k')
26.	
27.	# Decision boundary hyperplane
28.	db1 = hyperplane(hyp_x_min, svm.w, svm.b, 0)
29.	db2 = hyperplane(hyp_x_max, svm.w, svm.b, 0)
30.	plt.plot([hyp_x_min, hyp_x_max], [db1, db2], 'y--')
31.	
32.	plot_decision_boundary(svm)
33.	plt.show()

The `hyperplane` function calculates the y-values of the hyperplane for given x-values, weight vector `w`, bias `b`, and a value `v` which is either 1, -1, or 0. These values correspond to the positive support vector hyperplane, the negative support vector hyperplane, and the decision boundary, respectively.

The efficiency of the prediction process is notable. Once the model parameters  $w$  and  $b$  are determined, predictions can be made very quickly. For instance, predicting the class of a point  $(5, 2)$  or  $(5, -2)$  is instantaneous due to the simplicity of the computation:

1.	print(svm.predict([5, 2]))
2.	print(svm.predict([5, -2]))



Finally, it is essential to verify the values of the support vectors. These vectors lie close to the decision boundary and play a important role in defining it. The support vectors can be identified by examining the values of the margin constraints, which should be close to 1.

The SVM prediction phase involves classifying new data points based on the learned model parameters and visualizing the decision boundary to understand the classification results. This approach ensures a clear understanding of how the SVM model operates and its effectiveness in separating different classes.

In the context of optimizing a Support Vector Machine (SVM) from scratch using Python, it is critical to understand the importance of the optimization steps and the role of parameters such as the B range multiple. The optimization process involves iteratively refining the decision boundary to ensure that the SVM performs accurately and efficiently.

When iterating through the data to optimize the SVM, it is essential to check whether both classes have values close to the boundary. If one of the classes does not have a value less than a specific threshold, such as 1.02 or 1.001, further optimization is required. This involves taking additional steps to refine the decision boundary. It is evident that skipping these final optimization steps can significantly degrade the model's performance, as seen through the misclassification of points and a visibly shifted decision boundary.

The B range multiple is another important parameter that affects the optimization process. Setting the B range multiple to a lower value, such as one, can significantly speed up the initial optimization steps while maintaining reasonable accuracy. However, achieving a perfect classification may require a higher B range multiple, which increases computational expense. It is important to balance the B range multiple to achieve an optimal trade-off between accuracy and computational efficiency.

The optimization process can be implemented programmatically by checking if both classes have values close to the boundary in each iteration. If not, further tweaking is required. This can be achieved using a for loop to iterate through the possible values and checking for proximity to the boundary. Here is an example of how this check can be implemented in Python:

1.	for value in range(start, stop, step):
2.	if not (class1_value < threshold and class2_value < threshold):
3.	# Further optimization needed
4.	continue
5.	# Proceed with the current optimization step

Parallelization is another consideration in optimizing SVM. While certain steps in the optimization process cannot be threaded due to dependencies on prior knowledge, other parts of the code can be parallelized to improve efficiency. For example, saving results to a dictionary can be done in parallel without any issues. This parallelization can significantly speed up the optimization process.

Additionally, the concept of step sizes can be applied to the parameter B to find the maximal B efficiently. By making small changes to the range of B, significant gains in optimization speed can be achieved. This approach is analogous to the step sizes used in the primary optimization process.

Optimizing an SVM from scratch involves careful consideration of parameters such as the B range multiple and the application of iterative refinement steps. Parallelization and efficient parameter tuning can significantly enhance the optimization process, leading to a well-performing SVM model.

In machine learning, one of the critical scenarios involves handling data that is not linearly separable. Support Vector Machines (SVM) are a powerful tool for classification tasks, particularly when dealing with such complex datasets. The fundamental idea behind SVM is to find the hyperplane that best separates the classes in the feature space. This hyperplane is chosen to maximize the margin, which is the distance between the hyperplane and the nearest data points from each class, known as support vectors.

To implement SVM from scratch using Python, one must understand the optimization problem that SVM solves. The objective is to minimize the following function:

$$L(w, b, \xi) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i$$

subject to the constraints:

$$y_i(w \cdot x_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0$$

Here,  $w$  is the weight vector,  $b$  is the bias,  $\xi_i$  are slack variables for handling misclassifications,  $C$  is the regularization parameter,  $y_i$  are the class labels, and  $x_i$  are the feature vectors. The first term,  $\frac{1}{2} \|w\|^2$ , aims to maximize the margin, while the second term,  $C \sum_{i=1}^n \xi_i$ , penalizes misclassifications.

In Python, this can be implemented using libraries such as NumPy for numerical operations and optimization libraries like SciPy. Below is a simplified pseudocode to illustrate the process of implementing SVM:

1.	import numpy as np
2.	from scipy.optimize import minimize
3.	
4.	def svm_loss(w, X, y, C):
5.	n = X.shape[0]
6.	distances = 1 - y * (np.dot(X, w[:-1]) + w[-1])
7.	distances[distances < 0] = 0 # equivalent to max(0, distance)
8.	hinge_loss = C * np.sum(distances)
9.	return 0.5 * np.dot(w[:-1], w[:-1]) + hinge_loss
10.	
11.	def fit_svm(X, y, C):
12.	w_init = np.zeros(X.shape[1] + 1)
13.	opt_result = minimize(fun=svm_loss, x0=w_init, args=(X, y, C), method='L-BFGS-B')
14.	return opt_result.x[:-1], opt_result.x[-1]
15.	
16.	# Example usage
17.	X = np.array([[2, 3], [3, 3], [4, 5], [1, 1], [2, 1]])
18.	y = np.array([1, 1, 1, -1, -1])
19.	C = 1.0
20.	w, b = fit_svm(X, y, C)

In this pseudocode, `svm_loss` calculates the loss function incorporating the hinge loss and the regularization term. The `fit_svm` function uses the `minimize` function from SciPy's optimization library to find the optimal weight vector  $w$  and bias  $b$ .

Understanding these concepts and the mathematical formulation behind SVM is important for effectively applying this algorithm to classification problems. Mastery of the implementation details ensures that one can customize and extend the SVM algorithm to suit specific needs and datasets.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: KERNELS INTRODUCTION**

In this tutorial, we will continue our discussion on Support Vector Machines (SVM) in the context of Artificial Intelligence and Machine Learning with Python. Specifically, we will explore the concept of kernels and their role in SVM.

Up until now, we have been working with linearly separable data. However, in real-world scenarios, linear separability is often not achievable. To address this, we can take a different perspective by adding a new dimension to our feature set. By introducing a new dimension, we can potentially achieve linear separability.

To illustrate this, let's consider a two-dimensional feature set with two classes: plus and minus. In this case, it is not possible to determine the support vectors or find the best separating hyperplane. To overcome this, we can add a third dimension using a function, such as  $x_3 = x_1 * x_2$ . This additional dimension allows us to visualize the data in a different way.

However, when dealing with high-dimensional data, such as image or video analysis, adding dimensions becomes impractical. Increasing the dimensionality of the data by even a small amount can significantly impact the training process of SVM, which relies on quadratic programming and optimization. Therefore, multiplying the dataset by 1.5X or more is not ideal, as it weakens the algorithm's performance.

Fortunately, there is an alternative approach that allows us to perform calculations in plausibly infinite dimensions without actually visiting those dimensions or incurring additional processing costs. This is where kernels come into play.

Kernels are similarity functions that take two inputs and output their similarity. They are not exclusive to SVM but are commonly associated with it. Kernels can be used to augment or enhance SVM by transforming nonlinear data into a higher-dimensional space, where linear separability can be achieved. This transformation allows us to work with nonlinear data without explicitly increasing the dimensionality of the dataset.

The value of kernels lies in their ability to handle complex data without the need for excessive dimensionality expansion. In the previous example, we discussed the challenges of adding dimensions to achieve linear separability. If we were to add multiple dimensions, the computational complexity would increase exponentially. Kernels provide a more efficient solution by performing calculations in higher dimensions without explicitly visiting those dimensions.

It is important to note that kernels are based on inner product operations. This means that they leverage the similarities between data points to determine separability. By applying kernels, we can transform the data into a space where linear separability is possible, even in cases where the original data is nonlinear.

Kernels are a powerful tool in SVM that allow us to handle nonlinear data by transforming it into a higher-dimensional space. By leveraging the inner product operations, kernels enable us to achieve linear separability without explicitly increasing the dimensionality of the dataset. This approach offers a more efficient and effective solution for handling complex data in machine learning.

An inner product and a dot product are essentially the same thing. In Python, if you create a couple of vectors and use the "dot" function or the "inner" function from the numpy library, you will find that they give the exact same results. Some people may use the term "dot" while others may use "inner," but they are referring to the same operation. In the context of writing a paper, it is generally recommended to use the term "inner product." However, it is always a good idea to verify this information independently.

Now, let's discuss the use of kernels in support vector machines (SVMs). To determine if we can use a kernel, we need to establish if we can use an inner product. This is because using a kernel involves transforming our data into a new dimensional space. Up until now, we have been working in an "X" space, where our feature sets are denoted as  $X_1$ ,  $X_2$ , and so on, and we already have a value for "Y," which represents the class.

The next logical step would be to introduce a "Z" space. The question then arises: can we interchange "X" and

"Z"? Fundamentally, we can, but we need to consider if every interaction with the "X" space in our optimization algorithm and support vector involves a dot product or an inner product.

To answer this question, let's start at the end and work our way backwards. When we have an unknown feature set, denoted as "X," how do we determine its classification? We use the equation  $Y = \text{sign}(WX + B)$ . Here, "WX" represents the dot product between "W" and "X," and "B" is a scalar value. Since the result of "WX" is a scalar, it doesn't matter if "X" is in five dimensions or fifty dimensions.

Moving on to the constraints, there are two major constraints to consider. The first constraint involves the requirement that  $Y_{\text{sub } I} \text{ multiplied by } X_{\text{sub } I} \text{ dot } W + B - 1$  must be greater than or equal to 0. Again, we see that the interaction between  $X_{\text{sub } I}$  and "W" is a dot product. Therefore, we can replace  $X_{\text{sub } I}$  with a theoretical  $Z_{\text{sub } I}$  without any issues.

The second constraint involves finding values for "alpha" in the quadratic programming equation. This equation eventually gives us the sum over  $\alpha_{\text{sub } I}$  multiplied by  $Y_{\text{sub } I}$  multiplied by  $X_{\text{sub } I}$ . Once again, we observe that the interaction between  $X_{\text{sub } I}$  and  $\alpha_{\text{sub } I}$  is a dot product.

If we go back to the longer equation we initially wrote, we can see that all interactions involve a dot product. This confirms that every interaction in our support vector machine algorithm is a dot product.

The inner product and the dot product are essentially the same thing. In the context of support vector machines, we can use a kernel if we can use an inner product. Interchanging "X" and "Z" does not affect the classification algorithm, and all interactions in the algorithm involve a dot product.

In this tutorial, we will discuss the concept of kernels in the context of support vector machines (SVM) and machine learning. A kernel is essentially a similarity function that takes two inputs and outputs their similarity using the inner product. It is important to note that kernels are not unique to SVMs, but they are commonly used in this context.

The main advantage of using kernels is that they allow us to transform our feature space without incurring any additional processing cost. This means that we can effectively go from a lower-dimensional feature space to a higher-dimensional one without actually paying the price of computing the additional dimensions. While we may not always go to infinite dimensions, the ability to go to higher dimensions without increasing computation cost is a valuable feature of kernels.

To better understand the concept of kernels, let's consider an example where we have a feature set and want to transform it using a kernel. By visualizing the transformation, we can gain insight into how the kernel works. Additionally, we will work out the mathematical details to demonstrate why kernels are advantageous.

In the next tutorial, we will apply a kernel and work through the math by hand to gain a deeper understanding of how kernels function. We will specifically focus on the polynomial kernel and demonstrate its application. While there are some kernels that may be difficult to conceptualize, we will at least cover one of them to provide a general idea of their functionality.

Finally, we will move on to implementing kernels in Python. We will showcase an example of a kernel at work, highlighting its practical application in machine learning.

A kernel is a similarity function that can be used to transform feature spaces. By utilizing kernels, we can effectively increase the dimensionality of our feature space without incurring additional processing costs. Kernels are not specific to SVMs and can be applied to various machine learning algorithms or even used to create our own algorithms.

If you have any questions or comments, please feel free to leave them below. Thank you for watching, and we appreciate your support and subscriptions.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: REASONS FOR KERNELS**

In machine learning, support vector machines (SVM) are commonly used for handling non-linearly separable data. To achieve this, SVMs utilize the concept of kernels. Kernels allow us to transform the feature set into a higher-dimensional space where the data becomes linearly separable. In this didactic material, we will explore the reasons behind using kernels in SVMs.

When working with kernels, we typically denote the feature set as  $X$  and the new space as  $Z$ . The kernel function, denoted as  $K(X, X')$ , is responsible for mapping the feature set to the new space. It calculates the dot product between the transformed vectors  $Z$  and  $Z'$ . The dot product represents the similarity between the vectors and serves as the kernel.

The transformation from  $X$  to  $Z$  can be seen as applying a function to the feature set.  $Z$  is a function of  $X$ , and  $Z'$  is a function of  $X'$ . When performing a kernel operation, we convert  $X$  to  $Z$  and  $X'$  to  $Z'$ . Then, we calculate the dot product between  $Z$  and  $Z'$ . It is important to note that the functions applied to  $X$  and  $X'$  must be the same for the kernel to be valid. This symmetry ensures consistency in the transformation process.

In some equations, the kernel function may be denoted as  $\phi$  or  $\psi$ . For example, in the equation  $y = w * \phi(X) + b$ , the use of  $\phi$  indicates the utilization of a kernel instead of the traditional SVM calculation.

The dot product or inner product of  $Z$  and  $Z'$  yields a scalar value. This scalar value is all we need from the  $Z$  space. The question arises: can we calculate the inner product without explicitly knowing the  $Z$  space? The answer lies in the concept of feature mapping.

Let's consider a simple feature set in two dimensions,  $X_1$  and  $X_2$ . Suppose we want to transform this feature set to a second-order polynomial space. To achieve this, we start with a vector consisting of 1,  $X_1$ , and  $X_2$ . Then, we add the second-order terms, which are  $X_1^2$ ,  $X_2^2$ , and  $X_1 * X_2$ . This transformation results in a six-dimensional vector.

Kernels play an important role in SVMs for handling non-linearly separable data. They allow us to transform the feature set into a higher-dimensional space, where the data becomes linearly separable. The kernel function calculates the dot product between the transformed vectors, representing the similarity between them. By utilizing kernels, we can effectively classify complex data patterns.

Support Vector Machines (SVM) is a popular machine learning algorithm used for classification and regression tasks. In SVM, we often encounter the concept of kernels, which allow us to transform the data into a higher-dimensional space to make it easier to separate and classify. One commonly used kernel is the polynomial kernel.

To understand the polynomial kernel, let's first introduce the concept of the  $z$  space. In SVM, we start with a set of input vectors, denoted as  $X$ , where each vector is represented by its components  $x_1$  and  $x_2$ . The  $z$  space is a new space that we create by applying a second-order polynomial transformation to the input vectors. In the  $z$  space, each vector is represented as a new vector  $Z$ , which consists of the original components  $x_1$  and  $x_2$ , as well as additional terms obtained by squaring  $x_1$ , squaring  $x_2$ , and multiplying  $x_1$  and  $x_2$ .

The polynomial kernel, denoted as  $K(X, X')$ , is a function that calculates the dot product of two vectors in the  $z$  space, where  $X$  and  $X'$  are input vectors. The dot product of two vectors is obtained by multiplying the corresponding components of the vectors and summing them. In the case of the polynomial kernel, the dot product is calculated as follows:

$$K(X, X') = (1 + X \cdot X')^P$$

Here,  $\cdot$  represents the dot product operation, and  $P$  is the degree of the polynomial. In our previous example, we used a second-order polynomial, so  $P$  was equal to 2. However, you can choose any value for  $P$  depending on the complexity of the problem and the dimensionality of the data.

Using the polynomial kernel, we can avoid explicitly transforming the data into the  $z$  space and calculating the dot product. Instead, we can directly compute the kernel function using the input vectors  $X$  and  $X'$ . This saves computational resources and eliminates the need to deal with high-dimensional data.

To summarize, the polynomial kernel is a powerful tool in SVM that allows us to efficiently perform calculations in a higher-dimensional space without explicitly transforming the data. By choosing an appropriate degree for the polynomial, we can effectively separate and classify complex datasets.

Support Vector Machines (SVM) are powerful machine learning algorithms used for classification and regression tasks. In this tutorial, we will discuss the concept of kernels in SVM and specifically focus on two types of kernels: polynomial and radial basis function (RBF).

Polynomial kernels are used to transform the input data into a higher-dimensional feature space. This transformation allows SVM to find a linear decision boundary in this new space, even if the original data is not linearly separable. The polynomial kernel function is defined as  $K(X, X') = (1 + X \cdot X')^P$ , where  $X$  and  $X'$  are input vectors, and  $P$  is the degree of the polynomial. By increasing the degree of the polynomial, we can capture more complex patterns in the data. However, as the degree and the number of input features ( $n$ ) increase, the computational complexity also increases. This can make the equation more challenging to solve manually, but with the help of calculators or programming languages like Python, the increase in difficulty is not significant.

On the other hand, RBF kernels are used to transform the input data into an infinite-dimensional feature space. The RBF kernel function is defined as  $K(X, X') = \exp(-\gamma \|X - X'\|^2)$ , where  $X$  and  $X'$  are input vectors, and  $\gamma$  is a parameter that controls the smoothness of the decision boundary. The RBF kernel can capture more complex and non-linear patterns in the data. However, as the dimensionality increases towards infinity, it becomes harder to conceptualize and compute. In practice, going to infinite dimensions is not necessary, and the RBF kernel can handle most cases effectively.

Using a kernel to force data into linear separability can sometimes lead to overfitting. Overfitting occurs when the model learns the noise and outliers in the data, resulting in poor generalization to unseen data. In the next tutorial, we will discuss how to identify if overfitting has occurred and strategies to avoid it. It is essential to be cautious when using the RBF kernel, as it can also lead to overfitting in certain cases. However, in most scenarios, the RBF kernel works well and is the default kernel used in libraries like scikit-learn.

Kernels are an integral part of SVM, allowing us to transform data into higher-dimensional spaces and capture complex patterns. Polynomial kernels are suitable for cases where the data is not linearly separable, while RBF kernels can handle more complex and non-linear patterns. However, it is important to be aware of the potential for overfitting and to use appropriate strategies to avoid it.

Support Vector Machines (SVM) are powerful machine learning models used for classification and regression tasks. In SVM, the goal is to find the best hyperplane that separates the data points into different classes. However, sometimes the data may not be linearly separable, meaning that a straight line or hyperplane cannot accurately classify the data. This is where kernels come into play.

Kernels are a key component of SVM that allow us to transform the data into a higher-dimensional space where it becomes separable. By applying a kernel function to the input data, we can map it to a new feature space where the classes can be separated by a hyperplane. This process is known as the kernel trick.

There are different types of kernels that can be used in SVM, such as linear, polynomial, radial basis function (RBF), and sigmoid. The choice of kernel depends on the nature of the data and the problem at hand. Each kernel has its own characteristics and mathematical properties.

- Linear Kernel: The linear kernel is the simplest and most commonly used kernel. It works well when the data is linearly separable. The decision boundary is a straight line or hyperplane.

- Polynomial Kernel: The polynomial kernel maps the data to a higher-dimensional space using polynomial functions. It is useful when the data has complex relationships and the decision boundary is non-linear.

- RBF Kernel: The radial basis function (RBF) kernel is a popular choice for SVM. It transforms the data into an infinite-dimensional space using Gaussian functions. It is effective when the data is not linearly separable and

the decision boundary is complex.

- Sigmoid Kernel: The sigmoid kernel maps the data to a higher-dimensional space using sigmoid functions. It is suitable for binary classification problems. However, it is less commonly used compared to other kernels.

Choosing the right kernel is important for achieving good classification performance. It requires understanding the characteristics of the data and experimenting with different kernels to find the best fit. It is important to note that SVM with kernels can be computationally expensive, especially when dealing with large datasets.

Kernels are an essential component of Support Vector Machines that allow us to handle non-linearly separable data. By transforming the data into a higher-dimensional space, we can find a hyperplane that accurately separates the classes. The choice of kernel depends on the nature of the data and the problem at hand. Experimentation and understanding the data are key to selecting the appropriate kernel.



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SOFT MARGIN SVM**

Welcome to this educational material on support vector machines (SVM) in the context of machine learning with Python. In the previous videos, we discussed the concept of SVM and its application to non-linearly separable feature sets using kernels. We explored how kernels can be used to transform data into a higher dimension and perform dot products efficiently.

We specifically focused on the polynomial kernel, where we manually calculated the second-order polynomial and observed that performing dot products in the higher-dimensional space was time-consuming. However, when we used the polynomial kernel to mimic the second-order polynomial, we found that the computation was much quicker and simpler.

We then introduced the radial basis function (RBF) kernel, which can translate data into seemingly infinite dimensions. We discussed the potential challenges that can arise when using the RBF kernel and how to determine if the data is linearly separable or if another kernel should be considered.

Now, let's consider the concept of soft margin support vector machines. Imagine we have a dataset that is not linearly separable. In this case, we cannot draw a straight line to separate the data. However, by applying the RBF kernel, we can obtain a decision boundary that approximates the separation.

The support vector hyperplane represents the boundary, and it passes through the data points that are closest to it, known as support vectors. In our example, only two feature sets are not support vectors, indicating potential overfitting.

To address this issue, we can draw an alternative decision boundary that is almost a straight line. Although this new boundary may have a few violations, it will result in fewer support vectors and less overfitting. The support vector hyperplane for this alternative boundary will have two support vectors on either side.

Overfitting is a statistical problem that arises when a model fits historical data too closely, leading to poor generalization. By reducing the number of support vectors, we can mitigate the risk of overfitting and improve the model's performance on unseen data.

Soft margin support vector machines provide a way to handle non-linearly separable data by using kernels. By carefully selecting the kernel and adjusting the decision boundary, we can find a balance between accuracy and generalization.

In machine learning, support vector machines (SVM) are powerful algorithms used for classification and regression tasks. In this didactic material, we will focus on soft margin SVMs and how they handle data that is not perfectly separable.

When working with real-world data, it is important to consider that the future data may differ from the past data. Therefore, fitting the past data perfectly may lead to errors when predicting future data. To evaluate the performance of a trained model, we can use training and testing data. If the accuracy on the testing data is low, we might question the effectiveness of our model or consider using a different kernel.

One approach to assess the model's performance is to check the number of support vectors. Support vectors are the data points closest to the decision boundary. If the number of support vectors is large compared to the total number of samples, it indicates potential overfitting. Typically, if the number of support vectors exceeds 10% of the total samples, there is a higher chance of overfitting.

If the accuracy is low and the percentage of support vectors is high, trying a different kernel might be beneficial. On the other hand, if the accuracy is low and the percentage of support vectors is low, it suggests that the data may not be suitable for the current model. In such cases, exploring other options is recommended.

In the context of soft margin SVMs, we introduce the concept of a separating hyperplane. A soft margin classifier allows for some degree of error in classification. This is different from a hard margin classifier, which



strictly separates the data. In real-world scenarios, where data is often not perfectly separable, a soft margin classifier is more appropriate.

To incorporate error in the classification, we introduce slack variables. Slack allows some leeway in the constraint equation of the SVM. The original constraint equation was  $y \text{ sub } i \text{ multiplied by } X \text{ sub } i \text{ dotted with } W \text{ plus } B \text{ greater than or equal to } 1$ . With slack, we modify the equation to  $y \text{ sub } i \text{ multiplied by } X \text{ sub } i \text{ dotted with } W \text{ plus } B \text{ greater than or equal to } 1 \text{ minus the slack variable}$ .

The slack variable must be greater than or equal to 0. A value of 0 corresponds to a hard margin classifier, while larger values introduce more flexibility in the classification. The total slack is the sum of the individual slack variables.

A soft margin SVM is a useful tool for handling data that is not linearly separable. By allowing for a degree of error in classification using slack variables, we can create more flexible models that better accommodate real-world data.

A soft margin support vector machine (SVM) is a type of machine learning algorithm used for classification tasks. It is particularly useful when dealing with datasets that are not linearly separable, meaning that the classes cannot be separated by a straight line.

In the traditional formulation of SVM, the goal is to minimize the magnitude of vector  $W$ , which represents the hyperplane that separates the classes. However, in the case of soft margin SVM, we introduce a new parameter called  $C$ , which allows for some slack or errors in the classification.

The objective function of a soft margin SVM is to minimize one-half of the magnitude of vector  $W$  squared, plus  $C$  times the sum of all the slacks. The slacks represent the errors or violations of the margin, which is the region between the hyperplane and the closest data points of each class.

The parameter  $C$  plays a important role in determining the trade-off between minimizing the magnitude of vector  $W$  and reducing the number of violations. By increasing the value of  $C$ , we penalize violations more heavily, leading to a narrower margin and potentially overfitting the data. Conversely, decreasing the value of  $C$  allows for more violations and a wider margin.

It is important to note that  $C$  and the slacks have no direct relation to the magnitude of vector  $W$ , except in the context of this optimization problem. Therefore,  $C$  can be seen as a way to adjust the importance of minimizing the slacks relative to minimizing the magnitude of vector  $W$ , based on the specific requirements or preferences of the problem at hand.

To determine the optimal value of  $C$ , it is common practice to experiment with different values and compare the performance of the soft margin SVM on the same dataset. For example, running the algorithm with a high value of  $C$ , such as a million, and then comparing it to a lower value of  $C$ , such as one, can provide insights into the impact of  $C$  on the accuracy of the classifier.

In practical implementations, such as in scikit-learn, the default value for  $C$  is often set to one. However, it is important to adjust this parameter based on the specific characteristics of the dataset and the desired trade-off between accuracy and margin width.

A soft margin SVM is a powerful tool for classification tasks when dealing with non-linearly separable datasets. By introducing the parameter  $C$ , it allows for a flexible balance between minimizing the magnitude of vector  $W$  and reducing violations of the margin. Proper selection of  $C$  is important to avoid overfitting or underfitting the data.

In this tutorial, we will discuss the concept of Support Vector Machines (SVM) in the context of Machine Learning with Python. Specifically, we will focus on Soft Margin SVM.

Support Vector Machines are a type of supervised learning algorithm that can be used for classification and regression tasks. They are particularly effective in cases where the data is not linearly separable. SVMs work by finding a hyperplane in a high-dimensional feature space that separates the data into different classes.

Soft Margin SVM is an extension of the basic SVM algorithm that allows for some misclassifications in order to achieve a better overall fit. This is done by introducing a slack variable that allows data points to be on the wrong side of the margin or even on the wrong side of the hyperplane. The objective of the algorithm is to minimize the sum of the slack variables while still maximizing the margin between the classes.

To implement Soft Margin SVM in Python, we will be using the scikit-learn library. Scikit-learn provides a comprehensive set of tools for machine learning, including support for SVMs. We will start by importing the necessary libraries and loading our dataset.

Next, we will preprocess the data by scaling the features to have zero mean and unit variance. This step is important to ensure that all features are on a similar scale and to prevent any particular feature from dominating the optimization process.

Once the data is preprocessed, we can create an instance of the SVM classifier and specify the desired parameters. These parameters include the type of kernel to use, the regularization parameter, and the tolerance for convergence.

The kernel function determines the shape of the decision boundary and can be linear, polynomial, or radial basis function (RBF). The regularization parameter, also known as C, controls the trade-off between maximizing the margin and minimizing the misclassifications. A smaller value of C allows for more misclassifications, while a larger value of C enforces a stricter margin.

After training the SVM classifier on the data, we can evaluate its performance by making predictions on a test set and comparing them to the true labels. We can also visualize the decision boundary and the support vectors, which are the data points that lie on or within the margin.

Soft Margin SVM is a powerful algorithm for classification tasks in Machine Learning. By allowing for some misclassifications, it can handle complex datasets that are not linearly separable. In Python, we can implement Soft Margin SVM using the scikit-learn library, which provides a user-friendly interface for training and evaluating SVM models.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SOFT MARGIN SVM AND KERNELS WITH CVXOPT**

In this tutorial, we will be exploring the application of kernels in support vector machines (SVM) using the CVXOPT library in Python. We will also discuss the concept of soft margin SVM and visualize the impact of kernels on SVM.

Before we begin, it is important to note that the code used in this tutorial is not original and has been sourced from Matthew Blonde Belle's GitHub. Additionally, Christopher Bishop's book on pattern recognition and machine learning provides valuable information related to the topic.

CVXOPT is a useful library for this tutorial as it allows us to directly observe the impact of kernels on the SVM model. However, it is worth mentioning that CVXOPT may not be commonly used in practice, especially when working with support vector machines. In such cases, libraries like LIBSVM are typically preferred.

The first step is to understand the quadratic programming solver used in CVXOPT. The solver minimizes the equation:  $\frac{1}{2} * X^T * P * X + Q^T * X$ , subject to certain constraints. These constraints include  $G(X) \leq H$  and  $A * X = B$ . The solver allows us to optimize the SVM model by adjusting the values of the variables.

To gain a better understanding of the solver, you can refer to the documentation provided in the description or the text-based tutorial. There is also a simple example of solving a quadratic programming problem available on the CVXOPT website.

Moving on, let's take a look at the code. We start by importing the necessary libraries, such as numpy and CVXOPT. The kernels used in the SVM model are defined, including the Gaussian kernel and the polynomial kernel. The Gaussian kernel calculates the exponential of the squared Euclidean distance between two data points, divided by twice the value of the parameter Sigma. On the other hand, the polynomial kernel is the dot product of the input vectors, raised to the power of p.

The SVM model is initialized using the SVM object. The initialization method sets the kernel and the penalty parameter C. If C is set to None, it indicates a hard margin SVM. To create a soft margin SVM, simply assign a value to C.

The fit method of the SVM object prepares the values required for the quadratic programming solver. It solves the problem and obtains the solution, including the alphas (Lagrange multipliers), support vectors, intercept (bias), and the projection.

The prediction method of the SVM object calculates the projection of a new data point and returns the sign of the projection.

To visualize the results, we import the pylab library. Linearly separable data is generated for demonstration purposes.

This tutorial provides an overview of using CVXOPT and kernels in support vector machines. It demonstrates the impact of kernels on the SVM model and explains the concept of soft margin SVM. Although CVXOPT may not be commonly used, it offers valuable insights into the inner workings of SVMs.

Support Vector Machines (SVM) are a type of machine learning algorithm used for classification tasks. In this tutorial, we will discuss the concept of soft margin SVM and kernels with CVXOPT.

SVMs are typically used for linearly separable data, where a hyperplane can be drawn to separate the different classes. However, in real-world scenarios, data is often not linearly separable. This is where soft margin SVM comes into play. Soft margin SVM allows for some misclassifications in order to find a more flexible decision boundary.

To handle non-linearly separable data, SVMs use kernels. Kernels transform the input data into a higher-dimensional feature space, where it may become linearly separable. One commonly used kernel is the Gaussian

kernel, also known as the radial basis function (RBF) kernel. Other kernels, such as the polynomial kernel, can also be used.

CVXOPT is a Python library that provides tools for convex optimization. It can be used to solve the optimization problem involved in training SVMs with soft margin and kernels.

To demonstrate these concepts, we will go through some code examples. We will start with a linearly separable dataset and train a hard margin SVM using the default linear kernel. The code will plot the support vectors, which are the data points closest to the decision boundary.

Next, we will move on to a non-linearly separable dataset and train a soft margin SVM. This time, we will use a polynomial kernel instead of the default Gaussian kernel. The code will show that even with overlapping data, the SVM can still find a reasonable decision boundary.

To visualize how SVMs work with kernels, we will transform a two-dimensional dataset into a three-dimensional space using a kernel. This will allow us to see how the data becomes separable in the higher-dimensional space and how the decision boundary is represented.

Finally, we will discuss the SVM parameters and their impact on the model's performance. We will also briefly touch on how SVMs can be extended to handle multi-class classification problems.

Please note that the code examples and explanations provided in this tutorial assume some prior knowledge of SVMs and Python programming. If you have any questions or need further clarification, feel free to leave a comment.

Support Vector Machines (SVM) are a powerful class of machine learning algorithms used for classification and regression tasks. In this didactic material, we will focus on Soft Margin SVM and Kernels with CVXOPT.

Soft Margin SVM is an extension of the original SVM algorithm that allows for some misclassification errors in order to find a better decision boundary. This is useful when dealing with non-linearly separable data. The goal of Soft Margin SVM is to find the decision boundary that maximizes the margin while allowing for a certain number of misclassified points.

To achieve this, Soft Margin SVM introduces a regularization parameter, often denoted as  $C$ . This parameter controls the trade-off between maximizing the margin and minimizing the misclassification errors. A smaller value of  $C$  allows for a wider margin but may result in more misclassified points, while a larger value of  $C$  focuses on minimizing misclassification errors at the expense of a narrower margin.

Kernels are an essential component of SVM algorithms. They allow us to transform the input data into a higher-dimensional feature space, where it may become linearly separable. This is known as the kernel trick. The use of kernels enables SVM to handle complex data that cannot be easily separated in the original feature space.

CVXOPT is a Python library that provides an optimization framework for convex problems. It offers a user-friendly interface for solving the optimization problem associated with SVM. By utilizing CVXOPT, we can efficiently train Soft Margin SVM models and find the optimal decision boundary.

To summarize, Soft Margin SVM extends the original SVM algorithm by allowing for misclassification errors. Kernels play a important role in transforming the input data into a higher-dimensional feature space, making it easier to find a linear decision boundary. CVXOPT is a Python library that facilitates the optimization process associated with SVM.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: SUPPORT VECTOR MACHINE****TOPIC: SVM PARAMETERS**

In machine learning, support vector machines (SVM) are widely used for classification tasks. However, SVM is primarily a binary classifier, meaning it can only separate two groups at a time. But what if we have more than two groups to classify into? In such cases, we can use two methodologies: One versus Rest (OVR) and One versus One (OVO).

With OVR, we classify each group against the rest of the data. Let's consider an example where we have three groups: ones, twos, and threes. We start by separating the ones from the rest of the data using a separating hyperplane. Similarly, we separate the twos from the rest and the threes from the rest. Each group will have its own separating hyperplane.

However, OVR has a weighting issue. Each separating hyperplane is imbalanced because the number of data points on each side may differ. This can make it challenging to determine the correct classification for a given data point.

Alternatively, we can use OVO, where we create a separating hyperplane between each pair of groups. For example, we have a hyperplane separating ones from twos, another separating ones from threes, and another separating twos from threes. This approach eliminates the weighting issue of OVR.

When classifying a data point using OVO, we determine which side of each hyperplane the data point falls on. By considering the positions relative to all the hyperplanes, we can determine the most likely classification for the data point.

In practice, OVR is often the default choice, but both methodologies have their advantages and disadvantages. OVR is simpler and more straightforward, while OVO eliminates the weighting issue. The choice between the two depends on the specific problem and dataset.

Support Vector Machines (SVM) are powerful machine learning algorithms that can be used for classification tasks. In this didactic material, we will explore the parameters of SVM and how they can be adjusted to improve the performance of the model.

One important parameter in SVM is  $C$ , which determines the trade-off between maximizing the margin and minimizing the classification errors. A smaller value of  $C$  allows for a larger margin but may result in more misclassifications, while a larger value of  $C$  reduces the margin but leads to fewer misclassifications. It is important to note that  $C$  controls the soft margin classifier, and if a hard margin classifier is desired,  $C$  can be increased or decreased accordingly.

Another parameter in SVM is the kernel function, which is responsible for transforming the input data into a higher-dimensional space. The default kernel function is the radial basis function (RBF), but other options such as polynomial, linear, sigmoid, or even custom kernels can be used. The degree parameter is specific to the polynomial kernel and determines the power to which the kernel is raised. The gamma parameter is associated with the RBF kernel and controls the influence of each training example. It is recommended to leave gamma as the default value, which is calculated based on the number of features.

The independent term in the kernel function is represented by the parameter  $0$ . By default, it is set to 0, but it can be adjusted if necessary. SVM also provides the option to estimate probabilities. While SVM does not inherently provide a degree of confidence like other algorithms, probability estimates can be obtained by implementing cross-validation. However, it is important to note that this process can be computationally expensive.

Shrinking is another feature of SVM that can be enabled or disabled using the shrinking heuristic parameter. By default, it is set to true, and it improves the computational efficiency of the algorithm. Shrinking involves identifying feature sets that can be ignored during optimization, as they are deemed to have minimal impact on the final results.

Finally, the tolerance parameter determines the convergence criterion for the optimization algorithm. It controls the stopping criteria for the training process. A smaller tolerance value ensures a more accurate solution but may increase the training time.

SVM is a versatile machine learning algorithm that can be used for classification tasks. By adjusting the parameters such as C, kernel function, degree, gamma, independent term, probability estimates, shrinking, and tolerance, the performance of the SVM model can be optimized for different datasets and requirements.

In machine learning, support vector machines (SVMs) are powerful algorithms used for classification and regression tasks. SVMs work by finding an optimal hyperplane that separates data points into different classes. To achieve this, SVMs use various parameters that can be tuned to improve their performance. In this didactic material, we will focus on the SVM parameters and how they affect the optimization process.

One important question in SVM optimization is how to determine when we have reached the optimal solution. In SVM, the optimization process aims to find the best values for the weights (W) and the bias term (B) that minimize the classification error. To check if we have reached optimization, we can compare the values of the decision function, which is given by the equation  $y \text{ sub } i \text{ times } X \text{ sub } i \text{ dotted with } W \text{ plus } B \text{ minus } 1$ , with a tolerance value. If both sides of the equation have a value that is either  $1 \text{ e to the negative}$  or  $1 \text{ to the basically}$  name (e.g., 0.001), we can consider that we have found the optimal numbers and move on.

There are several SVM parameters that can be adjusted to improve the performance of the algorithm. One such parameter is the cache size, which determines the size of the kernel. The kernel is a mathematical function used to transform the input data into a higher-dimensional space, where it can be more easily separated. If the dataset is large, limiting the cache size can help prevent memory issues.

Another parameter is the class weight. By default, all classes are weighted equally. However, in some cases, certain classes may need to be given more importance. This parameter allows for adjusting the weights of different classes to reflect their importance in the classification task.

The max iterations parameter specifies the maximum number of iterations the optimization process will run. Each iteration involves updating the weights and bias based on the training data. Setting a higher value for max iterations allows for more iterations, potentially leading to a more accurate solution. However, it also increases the computational time.

The decision function shape parameter determines the strategy for multi-class classification. The options include one-versus-one (1v1), one-versus-rest (1vR), or none. The default option is none, which means that the algorithm uses a binary classification approach. However, for multi-class problems, it is recommended to use the one-versus-rest strategy.

SVM parameters play a important role in optimizing the performance of support vector machines. By adjusting these parameters, we can fine-tune the algorithm to achieve better classification accuracy. It is important to experiment with different parameter values to find the optimal combination for each specific problem.

Support Vector Machines (SVM) is a popular machine learning algorithm used for classification and regression tasks. In this didactic material, we will focus on the parameters of SVM and their significance in model performance.

One important parameter to consider is the decision function. In older versions of SVM, the decision function 'ovo' (one-vs-one) was used by default. However, in recent versions, 'ovo' is being deprecated and the default behavior is expected to change to 'OVR' (one-vs-rest). It is important to be aware of this change when working with SVM.

Another parameter to consider is the random state. This parameter is used to set a random seed for reproducibility. It is particularly useful when performing probability estimation. However, it is important to note that using probability estimation can significantly increase the computational load, especially for large datasets. Therefore, it is advisable to use probability estimation judiciously, considering the size of the dataset.

SVM also provides various attributes that can be useful for analysis and visualization. For example, the number of support vectors can be checked to evaluate the performance of the model. If the number of support vectors

is close to the total number of samples, it indicates that the model is not generalizing well. Additionally, the locations of the support vectors can be obtained, as well as the values of the weight vector ( $W$ ) and the bias term ( $B$ ). These attributes can be helpful for visualizing the decision boundary and understanding the model's behavior.

When working with SVM, it is important to consider the decision function, random state, and various attributes provided by the model. Understanding these parameters and attributes can help improve the performance and interpretability of the SVM model.



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: CLUSTERING INTRODUCTION**

Clustering is a technique in unsupervised machine learning where the machine is given a dataset and tasked with finding groups or clusters within the data. There are two major forms of clustering: flat clustering and hierarchical clustering.

In flat clustering, the scientist specifies the number of clusters they want the machine to find. For example, they may instruct the machine to find two or three clusters. On the other hand, hierarchical clustering allows the machine to determine the number of clusters and the groups within them.

One commonly used algorithm for clustering is the k-means algorithm. It is a simple algorithm that requires the scientist to choose the number of clusters, denoted as  $K$ . The algorithm works by randomly selecting  $K$  centroids, which are the centers of the clusters. These centroids can be chosen randomly or by using other methods such as shuffling the data.

Once the centroids are chosen, the algorithm calculates the distance between each data point and the centroids. This distance is typically calculated using the Euclidean distance formula. Each data point is then classified as belonging to the centroid it is closest to.

To illustrate this, let's consider an example with  $K=3$ . The first three data points are chosen as the initial centroids. The algorithm then calculates the distance between each data point and the centroids. Based on these distances, each data point is assigned to the centroid it is closest to.

The process is repeated iteratively, with the centroids being recalculated based on the new assignments. This continues until the centroids no longer change significantly or a maximum number of iterations is reached.

It's important to note that the initial choice of centroids can affect the final clustering result, and running the algorithm multiple times with different initial centroids can lead to different outcomes.

Clustering is a technique in unsupervised machine learning that involves finding groups or clusters within a dataset. The k-means algorithm is a popular method for clustering, where the scientist specifies the number of clusters to be found. The algorithm iteratively assigns data points to centroids based on their distances, until convergence is achieved.

Clustering is a technique used in machine learning to group similar data points together. In this didactic material, we will focus on the k-means clustering algorithm, which is one of the most popular and widely used clustering algorithms.

The k-means algorithm works by iteratively assigning data points to clusters and updating the cluster centers until convergence is achieved. The number of clusters, denoted as ' $k$ ', needs to be specified beforehand.

The algorithm starts by randomly initializing ' $k$ ' cluster centers. Then, for each data point, the algorithm calculates the distance to each cluster center and assigns the point to the cluster with the nearest center. This process is repeated until all data points have been assigned to a cluster.

Once all data points have been assigned to clusters, the algorithm calculates the mean of the feature sets or data points within each cluster, resulting in new cluster centers. This step is important as it helps to find the center of each cluster.

The process of assigning data points to clusters and updating the cluster centers is repeated until convergence is achieved. Convergence occurs when the cluster centers no longer move significantly. At this point, the algorithm has found the clusters.

It is worth noting that the k-means algorithm aims to cluster data into relatively equal-sized groups. This can be a limitation when dealing with datasets that contain clusters of different sizes. The algorithm's adherence to Euclidean distance may result in improperly sized clusters.



To address this limitation, other clustering algorithms use different approaches, such as using different kernels. However, even with different kernels, the issue of clustering differently sized groups may persist.

Scaling is another consideration when using clustering algorithms. Each data point needs to be compared to all other points, which can be computationally expensive. However, once the algorithm is trained and the cluster centers are determined, new points can be classified based on those centroids without the need for further training.

In practice, clustering algorithms can be used in semi-supervised machine learning. After clustering, the groups can be translated into supervised machine learning, where the clusters serve as labels for classification tasks using other algorithms like support vector machines.

To apply the k-means algorithm, we can use scikit-learn, a popular machine learning library in Python. Scikit-learn provides a straightforward implementation of the k-means algorithm, making it easy to apply to real-world examples.

The k-means clustering algorithm is a widely used technique for grouping similar data points together. It iteratively assigns data points to clusters and updates cluster centers until convergence. However, it may struggle with clustering differently sized groups due to its adherence to Euclidean distance. Scaling is also a consideration when using clustering algorithms, but once trained, new points can be classified based on the determined centroids. Scikit-learn provides an easy-to-use implementation of the k-means algorithm for practical applications.

In this didactic material, we will introduce the concept of clustering in the context of artificial intelligence and machine learning using Python. Specifically, we will cover the k-means and mean shift clustering algorithms.

To begin, we need to import the necessary libraries. We will use the matplotlib.pyplot module for data visualization and the sklearn.cluster module for clustering algorithms. We will also import numpy as np for numerical operations. The code for importing these libraries is as follows:

1.	<code>import matplotlib.pyplot as plt</code>
2.	<code>from matplotlib import style</code>
3.	<code>style.use('ggplot')</code>
4.	<code>import numpy as np</code>
5.	<code>from sklearn.cluster import KMeans</code>

Next, we will define a set of starting values for our data. These values will be stored in a numpy array. For example, we can use the following values: [1, 1.5, 1.85, 8.88, 1.06, 9.11]. These values represent the features of our data points. It is important to note that the number of elements in this array can vary depending on the desired number of clusters.

Once we have defined our data, we can visualize it using the scatter plot function from matplotlib.pyplot. We will plot the first and second elements of the array as the x and y coordinates, respectively. We will set the size of the markers to 150 and the line width to 5. The code for visualizing the data is as follows:

1.	<code>plt.scatter(x[:,0], x[:,1], s=150, linewidths=5)</code>
2.	<code>plt.show()</code>

After visualizing the data, we can proceed with the clustering process. We will use the k-means algorithm for this purpose. To create a k-means classifier, we need to define the number of clusters, which in this case is 2. We can do this by initializing an instance of the KMeans class with the parameter `n_clusters` set to 2. The code for creating the classifier is as follows:

1.	<code>clf = KMeans(n_clusters=2)</code>
----	---

Once we have created the classifier, we can fit it to our data using the fit method. This will assign each data point to one of the clusters based on their similarity. We can access the cluster centers and labels using the `cluster_centers_` and `labels_` attributes, respectively. The code for fitting the classifier and accessing the attributes is as follows:

1.	<code>clf.fit(x)</code>
2.	<code>centroids = clf.cluster_centers_</code>
3.	<code>labels = clf.labels_</code>

To visualize the clusters, we can assign different colors to the data points based on their labels. We can create a list of colors and assign a color to each label. For example, we can use green for label 0 and red for label 1. We can then plot each data point with its corresponding color. The code for visualizing the clusters is as follows:

1.	<code>colors = ["g.", "r."]</code>
2.	<code>for i in range(len(x)):</code>
3.	<code>plt.plot(x[i][0], x[i][1], colors[labels[i]], markersize=10)</code>
4.	<code>plt.scatter(centroids[:,0], centroids[:,1], marker="x", s=150, linewidths=5)</code>
5.	<code>plt.show()</code>

In the above code, we iterate over each data point and plot it with the corresponding color. We also plot the cluster centers as "x" markers.

By running this code, we can observe the clusters formed by the k-means algorithm based on the given data. The data points belonging to each cluster will be marked with the assigned color, and the cluster centers will be marked with "x" markers.

This concludes our introduction to clustering using the k-means algorithm in Python. We have covered the basic steps of importing libraries, defining data, visualizing data, creating a k-means classifier, fitting the classifier, and visualizing the clusters.

In this material, we will discuss the topic of clustering in the context of artificial intelligence and machine learning with Python. Specifically, we will focus on two popular clustering algorithms: k-means and mean shift.

Clustering is a technique used in machine learning to group similar data points together based on their characteristics. It is an unsupervised learning method, meaning that it does not require labeled data to train the model.

The k-means algorithm is one of the simplest and most widely used clustering algorithms. It aims to partition the data into k distinct clusters, where each data point belongs to the cluster with the nearest mean. The algorithm works iteratively, starting with an initial set of k centroids and assigning each data point to the nearest centroid. The centroids are then updated based on the mean of the data points in each cluster. This process is repeated until convergence, when the centroids no longer change significantly.

The mean shift algorithm, on the other hand, is a non-parametric clustering algorithm that does not require specifying the number of clusters in advance. Instead, it iteratively shifts the centroids towards the densest regions of the data. The algorithm starts with an initial set of centroids and computes the mean shift vector for each data point, which indicates the direction towards the densest region. The centroids are then updated by shifting them in the direction of the mean shift vector. This process is repeated until convergence.

To demonstrate these clustering algorithms, we will use Python. We will first visualize the data points and their initial centroids using a scatter plot. The data points will be colored based on their assigned clusters. We will then apply the k-means algorithm with different numbers of clusters to observe the resulting clusters and centroids. Similarly, we will apply the mean shift algorithm to see how it clusters the data.

In the next tutorial, we will apply clustering to an actual dataset. For example, imagine you work for a company like Amazon, and you have a dataset containing information about users. You believe that certain characteristics in the dataset can indicate whether a user is likely to make a purchase or not. By applying clustering algorithms like k-means, you can check if the algorithm separates the users into distinct groups based on their likelihood of being buyers. Later on, you can explore hierarchical clustering to identify different levels of likelihood instead of just binary groups.

By applying clustering algorithms to real datasets, we can gain insights and confirm the validity of certain assumptions. This can be valuable in various domains, such as customer segmentation, anomaly detection, and pattern recognition.

We hope you found this introduction to clustering informative. If you have any questions or comments, please feel free to leave them below. Thank you for watching, and we appreciate your support and subscriptions.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: HANDLING NON-NUMERICAL DATA**

In this tutorial, we will be discussing clustering, specifically flat clustering with the k-means algorithm. We will be using the Titanic dataset for this tutorial, which contains data from the passengers on the Titanic. The dataset includes information such as the passenger's class, whether they survived or not, their name, sex, age, number of siblings or spouses on board, number of parents or children, ticket number, fare, cabin, embarkation point, lifeboat number, body identification number, and their home or destination.

Our goal is to find insights from this data and determine if there are any patterns or relationships that can help us understand the survival rate of the passengers. We will be using the k-means algorithm to separate the passengers into two groups and analyze the survival rate of each group. Additionally, we will assess the accuracy of the clustering in predicting survival or death.

To begin, we will import the necessary libraries for preprocessing and cross-validation from the scikit-learn package. We will also import the pandas library as "PE" for data manipulation. Next, we will read the Titanic dataset into a data frame using the "read\_excel" function from pandas.

Once we have the data frame, we can examine the first few rows of the dataset using the "head" function to get a sense of the information it contains. It is worth noting that while some of the columns contain numerical values, such as passenger class and survival status, other columns, like the passenger's name, are not numerical. However, we can explore whether non-numerical data, such as the name, may have any significance in determining survival using techniques like natural language processing.

This tutorial focuses on applying the k-means clustering algorithm to the Titanic dataset to analyze the survival rate of different groups of passengers. We will also explore the potential significance of non-numerical data, such as the passenger's name, in predicting survival.

### Handling Non-Numerical Data in Machine Learning with Python

In machine learning, it is important to have numerical data for training models. However, sometimes we encounter non-numerical data in our datasets. In this didactic material, we will explore how to handle non-numerical data using Python.

One common example of non-numerical data is categorical data, such as gender or destination. These types of data cannot be directly used in machine learning algorithms. Therefore, we need to convert them into numerical form.

To convert non-numerical data into numerical data, we can use a technique called label encoding. Label encoding assigns a unique numerical value to each category in a column. For example, if we have a column with two categories, "female" and "male," we can assign "female" as 0 and "male" as 1.

To perform label encoding in Python, we can use the pandas library. First, we need to extract the column containing the non-numerical data. Then, we can use the set() function to get the unique values in the column. Next, we assign a unique numerical ID to each category in the set. Finally, we replace the original non-numerical data with the encoded numerical values.

It is important to note that if there are a large number of categories in a column, label encoding may result in outliers. These outliers can cause issues in machine learning algorithms. Therefore, it is recommended to perform data preprocessing to handle outliers before applying label encoding.

In addition to label encoding, we may also encounter missing data in non-numerical columns. To handle missing data, we have two options: dropping the rows with missing data or filling in the missing values.

Dropping rows with missing data may lead to a loss of valuable information. Therefore, it is often preferable to fill in the missing values. In Python, we can use the fillna() function from the pandas library to fill missing values with a specified value, such as 0.

When dealing with non-numerical data in machine learning, we need to convert the data into numerical form using techniques like label encoding. We should also handle missing data by either dropping rows or filling in the missing values.

Please note that this didactic material focuses on the technical aspects of handling non-numerical data and does not consider the theoretical background or advanced techniques. For further exploration, it is recommended to consult additional resources or educational materials.

In this didactic material, we will discuss how to handle non-numerical data in the context of artificial intelligence and machine learning using Python. Specifically, we will focus on clustering techniques such as k-means and mean shift.

When working with data, it is common to encounter non-numerical or categorical variables. These variables represent qualitative attributes rather than numerical values. To perform clustering analysis on such data, we need to convert these non-numerical variables into numerical form.

To handle non-numerical data in a data frame, we can follow a step-by-step process. First, we iterate through each column in the data frame. For each column, we create an empty dictionary called "text\_digit\_vals". This dictionary will store the mapping between the unique non-repetitive values in the column and their corresponding numerical representation.

Next, we define a function called "convert\_to\_int\_val" that takes an index value as input and returns the numerical representation of the corresponding non-numerical value. This function utilizes the "text\_digit\_vals" dictionary to perform the conversion.

We then check the data type of each column in the data frame. If the data type is not "int64" or "float64", indicating that it is a non-numerical variable, we proceed with the conversion. We create a list called "column\_content" containing the values in the column and obtain the unique elements from this list using the "set" function.

We iterate through each unique element and assign a numerical value to it using the "text\_digit\_vals" dictionary. This ensures that each unique non-numerical value is mapped to a unique numerical representation. Finally, we update the values in the column of the data frame using the "map" function and convert them to integers.

It is important to note that this approach may not be the most efficient for large datasets. However, it provides a simple and effective way to handle non-numerical data for clustering analysis.

To summarize, handling non-numerical data involves converting non-numerical variables into numerical form. This process allows us to apply clustering techniques such as k-means and mean shift to the data. By mapping unique non-numerical values to unique numerical representations, we can effectively analyze and cluster categorical variables in machine learning tasks.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: K MEANS WITH TITANIC DATASET**

In this tutorial, we will explore the concept of clustering in machine learning using the k-means algorithm with the Titanic dataset. Clustering is a technique used to group similar data points together based on their characteristics. The k-means algorithm is a popular clustering algorithm that aims to partition the data into a predetermined number of clusters.

Before we begin, it is important to note that clustering is an unsupervised machine learning technique, meaning that we do not have labeled data to train our model. Instead, we will use the features of the Titanic dataset to classify the passengers into two groups: survivors and non-survivors.

To start, we need to preprocess the data. We will convert any non-numerical data into numerical form, as this is a requirement for the k-means algorithm. We will also drop the "survived" column from the dataset, as including it would be considered cheating. To do this, we will use the "as\_type" function to convert the data frame to float.

Next, we will define our target variable, "y," which will contain the "survived" column of the data frame. This will be used to compare the groups identified by the clustering algorithm.

Now, we can proceed with applying the k-means algorithm. We will initialize the classifier, "CLF," with the k-means algorithm and specify that we want to create two clusters. We will then fit the algorithm to our preprocessed data, "X."

After fitting the algorithm, we can compare the groups identified by the k-means algorithm with the "survived" column. We will compare all the features except for the passenger's name and body identification number, as these are not relevant for our analysis. We will also consider excluding the "boat" feature, as whether or not a person was on a lifeboat may impact their chance of survival.

To compare the groups, we will iterate through the labels assigned by the k-means algorithm. For each label, we can predict whether the passenger survived or not. We can either use the "predict" function or iterate through the labels directly, as it will yield the same result.

In this tutorial, we have explored the concept of clustering using the k-means algorithm with the Titanic dataset. We have preprocessed the data, applied the k-means algorithm, and compared the identified groups with the "survived" column. This analysis can provide insights into the factors that may impact a passenger's chance of survival.

In this tutorial, we will explore the concepts of clustering, specifically k-means and mean shift, using the Titanic dataset. We will implement these clustering algorithms using Python and discuss their applications in machine learning.

Before diving into the implementation, let's briefly explain what clustering is. Clustering is an unsupervised learning technique that aims to group similar data points together based on their features. It helps in identifying patterns and similarities within the data, which can be useful for various purposes such as customer segmentation, anomaly detection, and image recognition.

First, let's focus on k-means clustering. The k-means algorithm aims to partition the data into k clusters, where each data point belongs to the cluster with the nearest mean. The algorithm works as follows:

1. Initialize k centroids randomly.
2. Assign each data point to the nearest centroid.
3. Recalculate the centroids as the mean of the data points in each cluster.
4. Repeat steps 2 and 3 until convergence.

To implement k-means clustering with the Titanic dataset, we start by importing the necessary libraries and loading the dataset. We then preprocess the data by converting non-numeric values to numeric ones. Next, we scale the data using the `preprocessing.scale()` function to ensure that all features have the same importance

during clustering.

We can now proceed with the clustering process. We create an instance of the KMeans class from the sklearn.cluster module and specify the number of clusters (k) we want to form. We fit the model to the scaled data using the fit() method and obtain the predicted labels for each data point using the predict() method.

To evaluate the accuracy of our clustering model, we compare the predicted labels with the actual labels from the dataset. Since clustering is an unsupervised learning technique, we don't have access to the true labels. However, in this case, we can compare the predicted labels with the "Survived" column, which indicates whether a passenger survived or not.

We calculate the accuracy by counting the number of correct predictions and dividing it by the total number of data points. By adding the preprocessing step of scaling the data, we can observe a significant improvement in the accuracy of the clustering model.

In addition to k-means clustering, we also briefly discuss mean shift clustering. Mean shift is another popular clustering algorithm that aims to find the dense regions of data points by iteratively shifting the centroids towards the higher density regions. It does not require specifying the number of clusters in advance, making it more flexible.

In this tutorial, we explored the concepts of clustering, specifically k-means and mean shift, using the Titanic dataset. We implemented these clustering algorithms using Python and discussed their applications in machine learning. By preprocessing the data and scaling the features, we observed improved accuracy in the clustering model.

In the previous material, we explored the impact of different columns in the Titanic dataset on the accuracy of our machine learning classifier. We observed that some columns, such as "Ticket" and "Boat," did not significantly affect the accuracy of our classifier, while others, like "Sex," had a significant impact.

By removing the "Boat" column, we obtained an accuracy of around 69%. This suggests that the column does not provide valuable information for predicting survival rates. Similarly, when we removed the "Sex" column, the accuracy dropped to around 68%. This indicates that gender is an important factor in determining survival rates.

Interestingly, even without explicitly instructing the classifier to separate individuals into groups based on their likelihood of survival, it was able to achieve an accuracy of approximately 70%. The classifier based its predictions on other features, such as class, age, siblings, ticket number, fare, cabin, and embarkation details.

We also discussed the potential impact of hierarchical clustering on the Titanic dataset. While we did not consider it in this material, hierarchical clustering can be used to identify a spectrum of groups within the data, rather than just two distinct classes. This approach may provide more nuanced insights into the dataset.

In the next tutorial, we will build our own k-means classifier. However, we will not be using the Titanic dataset for this exercise. We will instead explore hierarchical clustering to uncover additional information from the dataset. This will allow us to identify more interesting patterns and relationships within the data.

Our analysis revealed that certain columns, such as "Sex," have a significant impact on the accuracy of our machine learning classifier. By considering features like class, age, siblings, ticket number, fare, cabin, and embarkation details, the classifier was able to accurately predict survival rates with an accuracy of approximately 70%. This demonstrates the effectiveness of the k-means algorithm in clustering individuals into groups based on their likelihood of survival.



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: CUSTOM K MEANS**

In this didactic material, we will discuss the concept of k-means clustering in the context of machine learning with Python. Clustering is a technique that involves grouping similar data points together based on certain criteria. K-means clustering is a specific type of clustering algorithm where the number of clusters, denoted as  $K$ , is predetermined. In this material, we will focus on building a custom version of the k-means algorithm.

To understand k-means clustering, let's first visualize a two-dimensional graph with some data points on it. The goal of k-means clustering is to separate the data set into  $K$  number of groups. To achieve this, we start by randomly selecting two data points as the initial centroids. We then measure the distances of every other data point to these centroids. Each data point is classified as belonging to the centroid it is closest to. We calculate the mean of the data points in each class and update the centroids accordingly. We repeat this process iteratively until the centroids stop moving, indicating that the data has been successfully clustered.

To begin building our custom version of k-means, we will use some code from a previous tutorial (part 34). If you don't have the code, you can find it in the text-based version of this tutorial or in the provided material. We will define a class called "K\_means" and initialize it with the values of  $K$ , tolerance, and maximum iterations. The tolerance represents the maximum allowed movement of the centroids, and the maximum iterations determine how many times the algorithm will run before stopping.

Next, we will define the "fit" method, which will be responsible for training our custom k-means algorithm. The method takes in the data as input and performs the necessary calculations to cluster the data. We will also define the "predict" method, which will allow us to make predictions on new data using the trained centroids.

It's important to note that, unlike other machine learning algorithms, with k-means clustering, we can predict on the same data that was used for training. This is because the prediction is based on the centroids and not the individual data points. However, in a classification scenario, this would be considered cheating.

By building our own custom version of the k-means algorithm, we gain a deeper understanding of how the algorithm works and can customize it to suit our specific needs. This can be particularly useful when working with datasets that have unique characteristics or when we want to experiment with different variations of the algorithm.

K-means clustering is a powerful technique for grouping data points into distinct clusters. By building our own custom version of the k-means algorithm, we can gain a better understanding of its inner workings and adapt it to our specific requirements.

In this didactic material, we will discuss the concept of custom k-means clustering in machine learning with Python. Clustering is a technique used to group similar data points together based on their characteristics. K-means is a popular clustering algorithm that aims to partition the data into  $k$  distinct clusters. However, in some cases, we may want to customize the k-means algorithm to suit our specific needs.

To begin, we need to define the custom k-means algorithm. The first step is to initialize the centroids, which are the representative points for each cluster. In our implementation, we will start with an empty dictionary to store the centroids. Next, we will iterate through the desired number of clusters, denoted by  $k$ , and assign the initial centroids as data points from the dataset.

Once we have initialized the centroids, we move on to the optimization process. This involves iterating through a specified number of iterations, which can be set using the `max_iterations` parameter. During each iteration, we will update the classifications of the data points based on their distances to the centroids.

To calculate the distances, we will use the numpy library's `linalg.norm` function. This function calculates the Euclidean distance between a feature set and each centroid. The distances will be stored in a list for each data point.

Next, we will assign each data point to the centroid with the minimum distance. This is done by finding the



index of the minimum distance in the distances list and assigning the corresponding centroid as the classification for that data point. We will update the classifications dictionary accordingly.

After completing all iterations, we will have the final classifications for each data point. These classifications represent the clusters to which the data points belong. We can use this information for further analysis or visualization of the clustering results.

It is important to note that the starting centroids may not be optimal for clustering. In some cases, it may be necessary to shuffle the dataset and repeat the process to achieve better results. However, in most cases, the custom k-means algorithm performs well without the need for additional optimization.

Custom k-means clustering is a powerful technique for grouping data points into clusters based on their characteristics. By customizing the k-means algorithm, we can tailor the clustering process to our specific needs. This approach allows us to gain insights and make informed decisions based on the patterns and similarities within the data.

In the previous material, we discussed the process of implementing the custom K-means clustering algorithm using Python. We covered the steps involved in creating centroids and classifying data points based on their proximity to these centroids. In this section, we will continue our discussion and explore the remaining steps of the algorithm.

After initializing the centroids, we need to iterate through the data points and classify them based on their proximity to the centroids. To do this, we empty out the classifications and redo the classification process every time the centroid changes. This ensures that the classification is up to date with the latest centroid positions.

Once the classification is complete, we calculate the average of the feature values for each class. This is done using the `np.average` function, which takes the average of all the classifications for a given centroid. The `axis=0` parameter specifies that the average should be calculated along the columns of the dataset.

By finding the mean of all the features for any given class, we are able to determine the centroid for that class. This centroid represents the average feature values for the data points belonging to that class. This step is important in updating the centroids and finding how much they have changed.

To compare the current centroids with the previous centroids, we set `pre_centroids` equal to `self.centroids`. This is necessary due to object inheritance. If we simply set `pre_centroids` equal to `self.centroids`, it would always be equal to the current centroids, which is not what we want. We need to compare the current centroids with the previous centroids to determine how much they have changed.

In the next iteration, we update the centroids based on the new classifications. We calculate the average of the feature values for each class and assign it as the new centroid for that class. This process is repeated until the centroids no longer change significantly, which is determined by a tolerance value.

It is important to note that the first iteration is slightly different. In the first iteration, we do not update the centroids based on the new classifications. This is done to show the initial step of the algorithm, where centroids are created and data points are assigned to the closest centroids.

In the next video, we will continue our discussion and likely be able to finish implementing the entire custom K-means clustering algorithm. We will also test the algorithm to see if it produces the desired results. The remaining steps, such as the predict function, are relatively quick to implement. We appreciate your support and encourage you to leave any questions, comments, or concerns below. Thank you for watching!

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: K MEANS FROM SCRATCH**

In this didactic material, we will discuss the topic of clustering in the context of machine learning with Python. Specifically, we will focus on the k-means algorithm and mean shift algorithm. Clustering is a technique used to group similar data points together based on their inherent characteristics. It is commonly used in various applications such as customer segmentation, image recognition, and anomaly detection.

The k-means algorithm is an iterative algorithm that aims to partition a given dataset into k clusters. The algorithm starts by randomly selecting k centroids, which act as the initial cluster centers. It then assigns each data point to the nearest centroid based on a distance metric, such as Euclidean distance. After the assignment, the algorithm recalculates the centroids by taking the mean of all the data points assigned to each cluster. This process is repeated until convergence, where the centroids no longer change significantly.

To implement the k-means algorithm from scratch, we need to define a class that encapsulates the necessary methods and attributes. One important method is the fit method, which takes the dataset as input and performs the clustering process. Inside this method, we initialize the centroids, assign data points to clusters, and update the centroids iteratively until convergence. We also define a predict method, which assigns new data points to the clusters based on the learned centroids.

In addition to the k-means algorithm, we will also briefly discuss the mean shift algorithm. Mean shift is a non-parametric clustering algorithm that does not require specifying the number of clusters beforehand. Instead, it iteratively shifts each data point towards the mean of its neighboring points until convergence. The resulting clusters are determined by the convergence points.

To visualize the results of our clustering algorithms, we can use the matplotlib library in Python. We can plot the data points and color-code them based on their assigned clusters. We can also plot the centroids as markers to indicate the cluster centers.

Clustering is a powerful technique in machine learning that allows us to group similar data points together. The k-means algorithm and mean shift algorithm are two popular clustering algorithms that can be implemented using Python. By understanding these algorithms and their implementation, we can effectively apply clustering in various real-world scenarios.

In this didactic material, we will discuss the topic of Artificial Intelligence - Machine Learning with Python, specifically focusing on Clustering, k-means, and mean shift algorithms. We will explore the concept of K means from scratch, without relying on pre-existing libraries or packages.

Clustering is a technique used in unsupervised machine learning to group similar data points together based on their characteristics. The k-means algorithm is one such clustering algorithm that aims to partition the data into k distinct clusters. Mean shift is another popular clustering algorithm that iteratively shifts the centroids of clusters towards the densest regions of data.

To understand K means from scratch, let's start by considering an example. We have a dataset consisting of unknown data points and a set of known centroids. The goal is to classify the unknown data points based on their proximity to the centroids.

In the code snippet, we can observe the following steps:

1. We initialize the known centroids.
2. For each unknown data point, we predict its classification by finding the closest centroid using the k-means algorithm.
3. We plot the unknown data points on a scatter plot, with different markers and sizes representing their classifications.
4. The centroid positions are not updated in this code snippet.

If we want to explore the impact of adding new data points to the dataset, we can modify the code by

appending the unknown data points to the original list. This will result in a change in the centroid positions, as the algorithm adjusts to the new data.

Furthermore, we can print the distances for optimization purposes. By monitoring the percent change in distances over iterations, we can observe how the algorithm converges towards an optimal solution.

To illustrate the process, we can create a live graph showing the iterations and movements of the centroids. This can help visualize the algorithm's progression and how the clusters evolve.

By adjusting parameters such as the maximum number of iterations, we can control the convergence speed of the algorithm. Smaller values may result in incomplete clustering, while larger values may increase the computational time.

In addition to the k-means algorithm, we briefly touch upon the meshing of the k-means classifier onto the Titanic dataset. We compare the results obtained using our implementation with those from the scikit-learn library. Interestingly, our implementation appears to compute the results faster than scikit-learn, although the exact reason for this difference is unknown.

The provided code demonstrates the process of handling non-numeric data in the context of the Titanic dataset. It showcases the steps involved in applying the k-means algorithm to such datasets.

This didactic material has covered the topic of Artificial Intelligence - Machine Learning with Python, specifically focusing on Clustering, k-means, and mean shift algorithms. We have explored the concept of K means from scratch and discussed its implementation steps. We have also briefly touched upon the comparison between our implementation and the scikit-learn library.

In the field of Artificial Intelligence, Machine Learning plays a important role in enabling computers to learn and make decisions without being explicitly programmed. One important aspect of Machine Learning is clustering, which involves grouping similar data points together based on certain characteristics. In this didactic material, we will focus on two popular clustering algorithms: k-means and mean shift, and we will explore how to implement k-means from scratch using Python.

Clustering algorithms are unsupervised learning techniques, meaning that they do not require labeled data for training. Instead, they aim to find patterns and structures within the data based on their inherent similarities. The k-means algorithm is a popular and widely used clustering technique that partitions the data into a predetermined number of clusters, where each data point belongs to the cluster with the closest mean.

To implement the k-means algorithm from scratch, we need to follow a step-by-step process. First, we initialize the centroids randomly or using some heuristic. Then, we assign each data point to the nearest centroid based on the Euclidean distance. After that, we update the centroids by calculating the mean of all the data points assigned to each cluster. We repeat these two steps until convergence, where the centroids no longer change significantly.

It is worth noting that there is no need to perform a train-test split when using clustering algorithms since the focus is on understanding how well the data is clustered as a whole, rather than making predictions. In the absence of labels, we can evaluate the clustering performance by analyzing the accuracy of the clusters. In our case, the accuracy ranges from 65% to 75%.

In the implementation of k-means, we observed that our custom implementation runs faster than the scikit-learn version. This may be due to the fact that the scikit-learn version imports additional functionalities that are not necessary for our specific task. However, further investigation is required to confirm this hypothesis.

The k-means algorithm provides a simple and intuitive way to cluster data, and implementing it from scratch allows us to gain a deeper understanding of its inner workings. By following the steps outlined in this didactic material, you should be able to create your own custom k-means classifier in Python.

In the next tutorial, we will shift our focus to hierarchical clustering, where the algorithm determines the optimal number of clusters automatically. To achieve this, we will utilize the mean shift algorithm. Stay tuned for an exciting exploration of hierarchical clustering in our upcoming tutorial.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: MEAN SHIFT INTRODUCTION**

Mean shift is a hierarchical clustering algorithm used in machine learning. Unlike the k-means clustering algorithm, mean shift does not require the user to specify the number of clusters beforehand. Instead, mean shift automatically determines the number of clusters and their locations.

To understand mean shift, let's consider a simple data set. In k-means, we randomly select K feature sets as cluster centers. However, in mean shift, every feature set is considered a cluster center. Each data point is assigned a bandwidth or radius, which determines the distance around the data point that is considered part of its cluster.

For example, let's focus on one data point. We assign a radius around this data point, and all other data points within this radius are considered part of the same cluster. The mean of these data points is calculated, and this becomes the new cluster center. A new bandwidth is then determined based on the new cluster center. This process is repeated until convergence, where the cluster center no longer moves.

It's important to note that this process is applied to every feature set or data point. Each cluster center goes through the same steps, updating its cluster center and bandwidth until convergence is reached.

The mean shift algorithm allows for the automatic determination of the number of clusters and their locations. As the algorithm progresses, cluster centers are refined, and data points are assigned to the appropriate clusters. When convergence is achieved, the algorithm is considered optimized, and the clusters are finalized.

Mean shift is a powerful clustering algorithm that can be used in various applications, such as image segmentation, object tracking, and anomaly detection. Its ability to automatically determine the number of clusters makes it a flexible and efficient tool in machine learning.

Mean shift is a clustering algorithm that aims to find the centers of clusters in a dataset. It works by iteratively shifting the cluster centers towards the regions of higher density. The algorithm starts with an initial set of cluster centers and then slowly moves them towards a point of convergence. Once the cluster centers stop moving, the algorithm is considered to have converged and the clustering process is complete.

In the case of mean shift, all cluster centers are moved towards the same point of convergence. This means that all the cluster centers will eventually converge to the same point. The algorithm determines the convergence point based on the density of the data points.

Sometimes, a simple radius and bandwidth may not be sufficient for convergence. In such cases, different levels of bandwidth can be used to determine the convergence point. Data points within each level of bandwidth can be assigned different weights, with higher weights indicating greater importance. This allows for a larger range of data points to be considered, while still penalizing for distance away from the cluster center.

To illustrate the concept, let's consider a dataset without any apparent clusters. We can define different levels of bandwidth and assign weights to the data points based on their proximity to the cluster centers. By applying the mean shift algorithm, we can determine the convergence point and identify clusters within the dataset.

In Python, we can implement mean shift using libraries such as scikit-learn. The algorithm starts by creating some starting centers, which are not used in the rest of the code or the mean shift algorithm. These centers are only used to generate the starting sample data. The cluster centers found by mean shift can then be compared to the randomly generated data to assess the accuracy of the algorithm.

To visualize the results, we can use 3D graphing with the help of libraries like matplotlib. By plotting the data points and the cluster centers, we can observe the clusters and their colors. We can also compare the cluster centers obtained from mean shift with the randomly generated data to evaluate the accuracy of the algorithm.

Mean shift can be a computationally intensive algorithm, especially for large datasets. It is recommended to run the algorithm on a powerful computer or a server to avoid performance issues.

Mean shift is a clustering algorithm that iteratively shifts the cluster centers towards the regions of higher density. It aims to find the convergence point where all cluster centers stop moving. By assigning different weights to data points based on their proximity to the cluster centers, mean shift can identify clusters within a dataset.

Clustering algorithms are commonly used in research and data structuring tasks, rather than for visualization purposes. In the next tutorial, we will discuss the application of the mean shift algorithm to the Titanic dataset. The goal is to determine the number of groups within the dataset. Initially, we may assume that there are two groups, but it is possible that there are more. For instance, the dataset already indicates the presence of three passenger classes: first, second, and third class. Thus, it is plausible that the true number of groups on the ship is three. Further exploration and analysis can be conducted from this point onwards. In the upcoming videos, we will consider these concepts. If you have any questions or concerns, please feel free to ask. Thank you for your support and interest in this material.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: MEAN SHIFT WITH TITANIC DATASET**

In this material, we will continue our discussion on the mean shift clustering algorithm. We will be using Python to implement the algorithm and analyze the Titanic dataset.

To start, we need to revisit the code from part 36. We will copy the code from fitting the k-means algorithm and paste it here. However, we need to make a few modifications. Instead of using k-means, we will be using mean shift clustering. Additionally, we will make a copy of the original data frame before dropping unnecessary columns.

The reason for making a copy of the data frame is that we want to analyze the clustered data after applying the mean shift algorithm. It is easier to interpret the results when the data is labeled with meaningful names instead of numerical values.

To implement mean shift clustering, we need to modify the code where the k-means algorithm was used. We remove the parameter for the number of clusters since mean shift does not require this parameter. With these modifications, we are ready to proceed.

After applying the mean shift algorithm, we obtain the labels for each data point and the cluster centers. We will add a new column to the original data frame called "cluster group". Initially, this column will be empty. We will populate it with the labels obtained from mean shift clustering.

To populate the "cluster group" column, we iterate through the labels and assign the corresponding value to each row in the original data frame. This is done by using the "iloc" function to reference the row and column index. The value from the labels array is assigned to the corresponding row in the "cluster group" column.

Now, let's discuss the power of mean shift clustering. Unlike k-means, mean shift clustering does not require us to specify the number of clusters in advance. It automatically determines the number of clusters based on the data. This flexibility allows mean shift to potentially discover more meaningful patterns in the data.

To illustrate this, we will calculate the survival rates for each cluster group. The survival rate will be calculated as the percentage of survivors within each cluster group. This information can provide insights into the characteristics of different groups identified by the mean shift algorithm.

We have implemented the mean shift clustering algorithm using Python and applied it to the Titanic dataset. We have labeled the data points with meaningful names and analyzed the resulting clusters. The flexibility of mean shift clustering allows it to potentially uncover interesting patterns in the data.

In this didactic material, we will explore the concept of clustering in machine learning using Python. Specifically, we will focus on the mean shift algorithm and its application to the Titanic dataset.

To begin, let's create a new temporary data frame that contains only the data points belonging to a specific cluster group. We can achieve this by using a conditional statement. The temporary data frame will be a subset of the original data frame where the cluster group is equal to a specified value, let's say zero.

Next, we will create another conditional data frame called "survival\_cluster". This data frame will contain only the data points from the temporary data frame where the "survived" column is equal to one. In other words, we are interested in the survival rate of the specific cluster group.

To calculate the survival rate, we will divide the length of the "survived" column in the "survival\_cluster" data frame by the length of the "temp\_DF" data frame. This will give us the proportion of survivors in the specific cluster group.

Finally, we will print the survival rate to the console. It is important to note that due to the slight randomness of the mean shift algorithm, the results may vary each time the code is run.

Moving on, let's analyze the results obtained from applying the mean shift algorithm to the Titanic dataset. We have three major cluster groups: group zero, group one, and group two. Group zero has a survival rate of 37%, group one has a survival rate of 84%, and group two has a survival rate of 10%.

Based on these results, we can infer that group zero is mostly comprised of second-class passengers, group one consists mainly of first-class passengers, and group two is predominantly made up of third-class passengers. It is worth noting that being in first class does not guarantee survival, but it may have increased the chances due to the proximity to lifeboats.

To further investigate, we can print the original data frame filtered by cluster group. For example, we can print the data points belonging to group one, which we have identified as the first-class passengers. Similarly, we can print the data points belonging to group two, which we assume to be the third-class passengers.

The mean shift algorithm applied to the Titanic dataset has allowed us to identify different cluster groups with varying survival rates. By analyzing the data points within each cluster, we have observed patterns suggesting that the survival rates are influenced by the passenger class. However, it is important to remember that these observations are based on a subset of the dataset and may not represent the entire population accurately.

In this tutorial, we will explore the concept of clustering in machine learning using Python. Specifically, we will focus on two clustering algorithms: k-means and mean shift. We will apply these algorithms to the Titanic dataset to gain insights into the passengers' characteristics.

Let's start by examining the data. The dataset consists of information for 1280 individuals. The mean value of the class variable is 2.3, indicating that the majority of passengers belong to class 2. However, there are also first-class and third-class passengers present in the dataset.

We can further analyze the data by looking at the average age of each group. The average age of the first-class passengers is 37, while the average age of the third-class passengers is 39. This suggests that the first-class passengers tend to be slightly older than the third-class passengers.

To gain a better understanding of the data, we can examine other variables, such as fare. For example, we can observe that the fare for group 2 (mostly third-class passengers) ranges from 29 to 69 British pounds. On the other hand, group 0 (mostly second-class passengers) shows a minimum fare of 0, indicating that someone might have boarded the Titanic for free. The maximum fare for group 0 is 263 pounds.

We can also analyze the fare for first-class passengers in group 1. The minimum fare paid by first-class passengers in this group is 247 pounds, while the maximum fare is 512 pounds. The mean fare for this group is 312 pounds. These findings suggest that the first-class passengers in this group belong to a more elite category.

Comparing the mean fare of group 1 (29 pounds) with the mean fare of group 2 (which is slightly higher), we can observe that the passengers in group 2 paid slightly more for their fare. This raises the question of whether the fare has any correlation with the survival rate.

To further investigate, we can analyze the survival rate of first-class passengers in group 0. The survival rate for first-class passengers in cluster 0 is 60%. This indicates that 60% of the first-class passengers in this group survived the Titanic disaster.

To gain a more comprehensive understanding, we can also analyze the survival rates of second-class and third-class passengers in this group, as well as in other clusters.

By applying clustering algorithms to the Titanic dataset, we were able to identify different groups of passengers based on their characteristics, such as class and fare. We also explored the survival rates within these groups, providing valuable insights into the passengers' experiences during the disaster.

In this tutorial, we will explore the concept of clustering in machine learning using Python. Specifically, we will focus on the mean shift algorithm and its application to the Titanic dataset.

Clustering is a technique used to group similar data points together based on their characteristics. Mean shift is one such clustering algorithm that iteratively shifts data points towards the mean of their neighborhood until



convergence is achieved.

Let's begin by running the mean shift algorithm on the Titanic dataset. The dataset contains information about passengers, including their ticket class, fare, and survival status. Our goal is to identify groups of passengers with similar characteristics.

After running the algorithm, we observe that it initially creates three groups. The survival rates for these groups are 0%, 37%, 84%, and 100%. Running the algorithm again reveals an additional group with a 100% survival rate.

Upon closer inspection, we find that the second group likely represents the bulk of the data, while the fourth group consists of passengers who paid a higher fare. It is interesting to note that the fare amount seems to have a significant impact on the survival rate.

In addition, we discover that gender does not play a significant role in determining the groups. Even when we remove variables like lifeboat information, the groups remain relatively consistent. This finding is intriguing and warrants further investigation.

The mean shift algorithm successfully identifies distinct groups within the Titanic dataset. By analyzing these groups, we can gain insights into factors that influence survival rates, such as fare amount. This dataset is commonly used for various analyses and serves as a good example of how clustering can be applied to group similar data points.

Feel free to explore the dataset further and uncover more interesting patterns. Clustering is a powerful technique that can be applied to a wide range of datasets, providing valuable insights into the underlying structure of the data.



**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: MEAN SHIFT FROM SCRATCH**

In this didactic material, we will be discussing the concept of mean shift clustering algorithm in the context of artificial intelligence and machine learning with Python. Mean shift clustering is a popular unsupervised learning technique used for grouping similar data points together. It is particularly useful when dealing with unstructured or unlabeled data.

To begin, we will need to understand the basic steps involved in the mean shift algorithm. The first step is to assign each data point as a cluster center. Next, we calculate the mean of all the data points within a certain radius, known as the bandwidth, of each cluster center. This mean value becomes the new cluster center. We repeat this process until convergence, where clusters stop moving or merge with each other.

Let's dive into the implementation of the mean shift algorithm in Python. We will start by initializing a class called "MeanShift" and defining an "init" method. In this method, we set the bandwidth value, which determines the radius for clustering. For simplicity, we will use a hard-coded bandwidth value, but in practice, it can be dynamically adjusted.

Next, we define a "fit" method, which takes in the data as a parameter. Inside this method, we create an empty dictionary called "centroids" to store the cluster centers. We then set the initial centroids by iterating through the data and assigning each data point a unique ID as the key in the dictionary, with the data point itself as the value.

Now, we enter an infinite loop using "while True". Inside the loop, we create an empty list called "new\_centroids" to store the newly calculated cluster centers. We then iterate through all the known centroids and create an empty list called "within\_bandwidth" to store the data points within the radius of each centroid.

Next, we iterate through the data and check if each feature set is within the radius of the current centroid. We use the Euclidean distance formula, implemented using the "numpy" library, to calculate the distance between the feature set and the centroid. If the distance is less than the bandwidth, we add the feature set to the "within\_bandwidth" list.

After iterating through the data, we have a list of all the feature sets within the bandwidth of the current centroid. We can now calculate the mean of these feature sets to obtain the new centroid. We repeat this process for all the known centroids.

Finally, we check for convergence by comparing the new centroids with the existing centroids. If the centroids have stopped moving or have merged with each other, we have achieved convergence, and the mean shift algorithm is complete.

This implementation provides a basic understanding of the mean shift algorithm and how it can be implemented from scratch using Python. Further improvements can be made by incorporating max iterations and tolerance values, as well as dynamically adjusting the bandwidth.

In this didactic material, we will be discussing the concept of Mean Shift clustering algorithm in the context of Artificial Intelligence and Machine Learning with Python. Mean Shift is a popular unsupervised learning algorithm used for clustering data points into groups based on their similarity. We will explore the implementation of Mean Shift from scratch, without using any pre-built libraries or functions.

To begin, let's clarify the difference between bandwidth and radius. Bandwidth refers to the range or spread of data points around a central point, while radius represents the distance from a central point to a data point. It is important to note that bandwidth encompasses the entire range, whereas radius is a specific distance.

In the implementation of Mean Shift, we start by calculating the mean of the centroid. The centroid represents the center point of a cluster. We iterate through each centroid and check if the norm (distance) between a data point and the centroid is less than the radius. If it is within the radius, we consider it to be within the bandwidth and append it to the feature set.

After appending the data points within the bandwidth to the feature set, we recalculate the mean of the centroid using the NP average function. This gives us the mean vector of all the vectors within the bandwidth. We then add this new centroid to a new centroids list.

Next, we obtain the unique elements from the new centroids list. We use the set function to get the unique tuples from the list and then convert it back to a sorted list. This step is important for convergence, as it removes duplicate centroids that are identical copies of each other.

To ensure that we can compare the previous centroids with the new centroids, we convert the previous centroids to a dictionary using the dict function. This allows us to save and modify the centroids without affecting the previous centroids.

Within the while loop, we define a new centroids dictionary and iterate through the unique elements. For each unique centroid, we convert it back to an array. We assume the centroids are optimized unless we find a reason why they are not. To check for movement, we compare the elements of the previous and new centroids arrays. If they are not equal, it indicates that there is movement and the centroids are not yet optimized.

This process continues until convergence is achieved, meaning that there is no further movement and the centroids are optimized.

Mean Shift is an unsupervised learning algorithm used for clustering data points based on their similarity. It involves calculating the mean of centroids, checking for data points within the bandwidth, and updating the centroids until convergence is achieved. By understanding the concepts and implementation of Mean Shift, we can effectively apply it to cluster data points in various machine learning tasks.

In this tutorial, we will continue our discussion on the topic of Artificial Intelligence and Machine Learning with Python. Specifically, we will focus on the concepts of clustering, k-means, and mean shift. In this section, we will explore mean shift from scratch.

To begin, we need to set the 'optimized' variable to False. This will allow us to enter the while loop. Inside the loop, we will iterate through the centroids and calculate the distances between each centroid and the points in the dataset. If a centroid has moved, we will update its position and set 'optimized' to False. This process will continue until all the centroids have converged.

To optimize the code and save processing time, we can add a check to break the loop if the centroids have not moved. This is done by using the 'if not optimized' statement. Additionally, we can break the for loop if 'optimized' is True, as there is no need to continue iterating.

Once the centroids have converged, we will reset their positions. This is done by setting 'self.centroids' to the final centroid positions.

Next, we will define the 'predict' function, although it will not be populated at this time. Following that, we will create an instance of the mean shift algorithm and fit it to the dataset. The centroids will be obtained from the 'centroids' attribute of the mean shift algorithm.

To visualize the results, we will scatter the data points and the centroids on a plot. The data points will be scattered using the 'scatter' function, while the centroids will be scattered using a different color and marker style.

After running the code, we can observe the cluster centers on the plot. The clustering algorithm appears to have worked as intended. However, it is important to note that the results can vary depending on the chosen radius value.

In the next tutorial, we will address the issue of determining a suitable radius automatically. We will also explore the concept of applying weights to different data points based on their proximity to the cluster center. This will help improve the accuracy of the clustering algorithm.

Thank you for following along with this tutorial. If you have any questions or comments, please feel free to leave

them below. Stay tuned for the next tutorial, where we will delve deeper into these topics.

**EITC/AI/MLP MACHINE LEARNING WITH PYTHON DIDACTIC MATERIALS****LESSON: CLUSTERING, K-MEANS AND MEAN SHIFT****TOPIC: MEAN SHIFT DYNAMIC BANDWIDTH**

In this tutorial, we will discuss the mean shift algorithm for clustering in the context of artificial intelligence and machine learning with Python. Specifically, we will focus on the concept of mean shift dynamic bandwidth.

Previously, we built our own custom mean shift algorithm, which worked well when the radius was set to 4 or around 4. However, if we increased or decreased the radius, the algorithm did not perform as intended. This limitation led us to address the issue of hard coding the radius value.

To overcome this limitation, we will introduce a dynamic bandwidth approach using weights. Instead of using a fixed radius, we will use a large radius and penalize points based on their distance from the current centroid. This approach will allow us to automatically handle finding the centroids correctly and make the code more dynamic.

To implement this, we will first set the default radius to None and define a radius norm step of 100. The radius norm step will determine the number of steps or bandwidths we will consider.

Next, we will modify the fitment function. If the radius is set to None, we will calculate the centroid of all the data points using the numpy average function. We will also calculate the magnitude from the origin to the data centroid using the numpy linalg.norm function.

With the average and the magnitude calculated, we can determine a suitable overall radius. We will divide the data norm by the radius norm step to obtain the new radius value.

In the while loop, before iterating through the feature set, we will define the weights. The weights will range from 0 to 99 and will be reversed using the negative one index. This means the weights will start from 99 and decrease to 0.

By incorporating these weights, we can assign higher values to data points closer to the centroid, effectively penalizing points further away from the centroid.

This mean shift dynamic bandwidth approach allows us to automatically handle finding the centroids correctly without the need for hard coding the radius. By penalizing points based on their distance from the centroid, we can achieve better results and make the code more dynamic.

Clustering is a fundamental technique in machine learning that allows us to group similar data points together. One commonly used clustering algorithm is k-means, which aims to partition the data into k distinct clusters. Another popular algorithm is mean shift, which identifies the densest regions of data points and assigns them as cluster centers. In this didactic material, we will focus on mean shift clustering and specifically discuss the concept of mean shift dynamic bandwidth.

Mean shift clustering involves iteratively shifting the cluster centers towards the densest regions of data points until convergence is reached. A important parameter in mean shift clustering is the bandwidth, which determines the size of the region considered for density estimation. In traditional mean shift clustering, a fixed bandwidth is used. However, in some cases, using a dynamic bandwidth can lead to better results.

To understand mean shift dynamic bandwidth, let's dive into the code implementation. In the code, we start by initializing the weights for each feature set in the data. These weights are calculated based on the distance between the feature set and the centroid. The closer the feature set is to the centroid, the higher the weight assigned to it.

To calculate the weight, we first compute the distance between the feature set and the centroid. If the distance is zero, we set it to a small value (0.01) to avoid division by zero. Next, we calculate the weight index by dividing the distance by the radius. The radius represents the distance within which we want to consider data points for density estimation.

If the weight index exceeds the maximum number of steps (`radius_norm_step`), we set it to the maximum value. This ensures that data points outside the desired radius are not considered. Finally, we square the weight index and multiply it by the feature set to obtain the weighted feature set.

It's important to note that this implementation may result in a large list of weighted feature sets, potentially increasing memory usage. To optimize this, one can consider alternative methods, such as calculating averages without creating a large array. However, for the purpose of showcasing the concept of weights, this simple implementation suffices.

Moving forward in the code, we accumulate the weighted feature sets in the bandwidth list using the `+=` operator. This operation allows us to add two lists element-wise. The resulting bandwidth list will be used to update the centroids.

To update the centroids, we compute the average of the feature sets in the bandwidth list using the `np.average` function. This gives us the new centroid positions. We also update the unique centroids by removing any duplicate entries.

In the subsequent part of the code, we initialize an empty list called `to_pop`. We then iterate over the unique centroids and for each centroid, we iterate over the unique centroids again. This nested loop structure is used to compare each centroid with every other centroid.

This comparison allows us to identify centroids that are close to each other and can be considered as a single cluster. If two centroids are close, we add them to the `to_pop` list. Finally, we remove the centroids in the `to_pop` list from the unique centroids.

Mean shift dynamic bandwidth is a technique that adjusts the bandwidth parameter in mean shift clustering based on the distance between data points and centroids. By assigning weights to feature sets and updating the bandwidth accordingly, we can achieve more accurate clustering results.

In this didactic material, we will discuss the concept of mean shift dynamic bandwidth in the context of artificial intelligence and machine learning with Python. Mean shift is a clustering algorithm that aims to find the modes or peaks of a density function, which can be interpreted as cluster centroids. By using a dynamic bandwidth, mean shift allows for the detection of clusters with varying densities.

To understand mean shift dynamic bandwidth, let's first recap the purpose of this algorithm. When working with large datasets, it is common to encounter centroids that are close to each other but not exactly identical. These centroids may differ by a small margin, which could be considered negligible. In such cases, it is unnecessary to have separate centroids for these similar data points. The goal of mean shift dynamic bandwidth is to identify and eliminate these redundant centroids.

To achieve this, we compare the distance between centroids using the `"norm"` function from the NumPy library. If the distance between two centroids is less than or equal to a specified radius, we consider them to be within one step of each other. In the context of our data set, which consists of a hundred steps, this means that if two centroids are within one original step, they need to be converged to the same centroid.

To implement this logic, we use a `"to_pop"` list to keep track of the redundant centroids. If two centroids are within the specified radius, we append the index of the redundant centroid to the `"to_pop"` list. After iterating through the data set, we remove the redundant centroids using the `"unique_start"` function. It is worth noting that modifying a list while iterating through it can lead to errors, so we use a try-except block to handle any potential issues.

Once we have removed the redundant centroids, we update the `"self.centroids"` list to reflect the changes. At this point, the major code changes for mean shift dynamic bandwidth are complete, and we can proceed to run the algorithm on our data set.

However, achieving perfect clustering may not always be possible, especially when dealing with small data sets. The presence of a limited number of data points can skew the centroids slightly. To address this, we can classify the clusters based on their cluster centers. By doing so, we can determine the success of the clustering algorithm even if the clusters are not perfect.

To implement the classification, we initialize an empty dictionary called "self.classifications". We then iterate through the range of the length of "self.centroids" and assign an empty list to each cluster index in "self.classifications". This allows us to classify the data points based on their proximity to the cluster centroids.

Finally, we calculate the distances between each data point and the cluster centroids using the "norm" function. We store these distances in a list of lists called "distances". Each sublist corresponds to a data point and contains the distances from that point to all the centroids. By comparing these distances, we can assign each data point to the closest centroid and update the "self.classifications" dictionary accordingly.

Mean shift dynamic bandwidth is a powerful algorithm for clustering data sets with varying densities. By eliminating redundant centroids and classifying the data points based on their proximity to the cluster centers, we can achieve effective clustering results. This algorithm can be implemented using Python and the NumPy library.

In the context of Artificial Intelligence and Machine Learning with Python, we will now discuss the topic of clustering, specifically focusing on k-means and mean shift algorithms. Clustering is a technique used to group data points into clusters based on their similarities. It is an unsupervised learning method commonly used in various applications such as customer segmentation, image segmentation, and anomaly detection.

The k-means algorithm is one of the most popular clustering algorithms. It aims to partition a dataset into k clusters, where k is a user-defined parameter. The algorithm works by iteratively assigning each data point to the nearest centroid and then updating the centroids based on the new assignments. This process continues until convergence, where the centroids no longer change significantly. The final result is a set of k clusters, each represented by its centroid.

The mean shift algorithm is another clustering algorithm that does not require the number of clusters to be specified in advance. Instead, it automatically determines the number of clusters based on the data. The algorithm starts by randomly selecting data points as initial centroids. It then iteratively shifts each centroid towards the region with a higher density of data points. This process continues until convergence, where the centroids no longer move significantly. The final result is a set of clusters, each represented by its centroid.

One important aspect of the mean shift algorithm is the concept of dynamic bandwidth. The bandwidth determines the size of the region around each data point that is considered when calculating the density. In the traditional mean shift algorithm, a fixed bandwidth is used. However, in the mean shift with dynamic bandwidth, the bandwidth is adaptively adjusted based on the local density of data points. This allows the algorithm to handle clusters of different sizes and densities more effectively.

To implement clustering algorithms in Python, we can use the scikit-learn library, which provides a wide range of machine learning algorithms and tools. The code snippet below demonstrates how to implement k-means and mean shift clustering using scikit-learn:

1.	from sklearn.cluster import KMeans, MeanShift
2.	import numpy as np
3.	
4.	# Generate sample data
5.	X, _ = make_blobs(n_samples=100, centers=3, n_features=2)
6.	
7.	# Perform k-means clustering
8.	kmeans = KMeans(n_clusters=3)
9.	kmeans.fit(X)
10.	kmeans_labels = kmeans.labels_
11.	
12.	# Perform mean shift clustering with dynamic bandwidth
13.	meanshift = MeanShift(bandwidth=None)
14.	meanshift.fit(X)
15.	meanshift_labels = meanshift.labels_
16.	
17.	# Visualize the clusters
18.	import matplotlib.pyplot as plt
19.	

20.	<code>plt.scatter(X[:, 0], X[:, 1], c=kmeans_labels)</code>
21.	<code>plt.title("K-means Clustering")</code>
22.	<code>plt.show()</code>
23.	
24.	<code>plt.scatter(X[:, 0], X[:, 1], c=meanshift_labels)</code>
25.	<code>plt.title("Mean Shift Clustering")</code>
26.	<code>plt.show()</code>

In the code snippet above, we first generate a sample dataset using the `make_blobs` function from scikit-learn. We then create instances of the `KMeans` and `MeanShift` classes and fit them to the data using the `fit` method. Finally, we visualize the clusters using scatter plots.

It is important to note that clustering algorithms may not always produce perfect results, especially when dealing with complex or overlapping data. It is often necessary to experiment with different parameter settings and preprocessing techniques to achieve the desired clustering outcome.

Clustering is a powerful technique in machine learning that allows us to group similar data points together. The k-means and mean shift algorithms are popular choices for clustering tasks. By understanding the concepts behind these algorithms and implementing them in Python using libraries like scikit-learn, we can effectively analyze and interpret complex datasets.

In this didactic material, we will discuss the concept of Mean Shift algorithm in the context of clustering in Machine Learning using Python. Mean Shift is a non-parametric clustering algorithm that aims to find the modes or dense regions of a dataset. Unlike other clustering algorithms such as k-means, Mean Shift does not require the number of clusters to be specified in advance.

Mean Shift works by randomly selecting data points as initial cluster centers and iteratively updating these centers based on the mean shift vector. The mean shift vector is calculated as the weighted average of the data points within a certain radius of each cluster center. The weights are determined by a kernel function, which assigns higher weights to data points closer to the cluster center.

The bandwidth or radius parameter of Mean Shift determines the size of the region used to calculate the mean shift vector. In the original Mean Shift algorithm, a fixed bandwidth is used. However, in practice, it is often challenging to select an appropriate fixed bandwidth that works well for all datasets. To address this issue, a dynamic bandwidth approach can be employed, where the bandwidth is adaptively adjusted based on the density of the data points.

In the transcript, the speaker runs the Mean Shift algorithm on a dataset and observes the resulting clusters. They note that the standard deviation of the clusters is significant, indicating the need to adjust the bandwidth parameter. The speaker also mentions the slow execution speed of their custom implementation compared to scikit-learn's implementation.

To improve the Mean Shift algorithm, the speaker suggests modifying the weights, step size, and squaring of the weights. They also propose using a weight dictionary to handle data points with different weights more efficiently. The speaker acknowledges that their implementation is not optimized and invites suggestions for improvement.

It is important to note that Mean Shift is a simple example of a clustering algorithm and may not be suitable for winning competitions. However, it provides an understanding of the algorithm's working principles and the challenges involved in parameter selection and optimization.

In the next material, we will transition to the topic of neural networks, which have gained popularity in recent years, particularly in the field of deep learning. Neural networks are initially simple but can become complex quickly. Stay tuned for our discussion on neural networks in the upcoming material.