

Task 6. Data Science example. Implement a sample machine learning program for a problem statement of your choice.

Explanation of the code:

This code demonstrates a basic example of using the scikit-learn library for machine learning tasks, specifically for linear regression. Here is a breakdown of the code line by line:

```
import pandas as pd
import numpy as np

from sklearn.datasets import fetch_california_housing
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
```

The code begins by importing the necessary libraries, including Pandas and NumPy for data manipulation, Scikit-learn for machine learning, and specifically the California housing dataset, linear regression model, and mean squared error metric.

Load the dataset

```
california_housing = fetch_california_housing(as_frame=True)
df = california_housing.frame
```

The `fetch_california_housing` function from Scikit-learn is used to load the California housing dataset. The `as_frame=True` argument specifies that the dataset should be returned as a Pandas DataFrame, which is assigned to the variable `california_housing`. The `frame` attribute is then used to extract the DataFrame object and assign it to the variable `df`.

Split the dataset into training and testing sets

```
X_train, X_test, y_train, y_test = train_test_split(df.drop('MedHouseVal', axis=1), df['MedHouseVal'],
test_size=0.2, random_state=42)
```

The `train_test_split` function from Scikit-learn is used to split the dataset into training and testing sets. The input data is split into two parts, one for training the model and another for testing it. Here, the `drop` function is used to remove the target variable, `MedHouseVal`, from the input data and assign it to `X_train` and `X_test`, while `y_train` and `y_test` are assigned to the target variable. The `test_size` argument specifies the proportion of the dataset to include in the testing set (in this case, 20%), and the `random_state` argument ensures that the data is split in the same way each time the code is run.

Train the model

```
lr = LinearRegression()
```

```
lr.fit(X_train, y_train)
```

A linear regression model is created using the `LinearRegression` function from Scikit-learn, and is assigned to the variable `lr`. The `fit` method is then called on the model object, passing in the training data as inputs.

```
# Evaluate the model
```

```
y_pred = lr.predict(X_test)
```

```
mse = mean_squared_error(y_test, y_pred)
```

```
print(f"Mean Squared Error: {mse}")
```

The `predict` method is used to generate predictions on the testing set using the trained model. The mean squared error metric is then calculated using the `mean_squared_error` function from Scikit-learn, which takes the true target values (`y_test`) and predicted target values (`y_pred`) as inputs. Finally, the mean squared error is printed to the console using an f-string. This metric provides a measure of how well the model is performing on unseen data, with lower values indicating better performance.

Output of the Code:

The output of this code is a single line of text that displays the mean squared error (MSE) of the linear regression model when applied to the testing data:

```
Mean Squared Error: 0.555891598695244
```

The MSE value indicates the average squared difference between the predicted target values (`y_pred`) and the actual target values (`y_test`) in the testing dataset. A lower MSE value indicates that the model is performing better in terms of accuracy, while a higher value indicates poorer performance. In this case, the MSE value is 0.5337, which suggests that the model has fairly good performance on the California housing dataset.

how these code can be helpful for Predicting House Prices with California Housing Dataset:

This code is an implementation of a linear regression model on the California Housing dataset, which contains information about various housing features in different districts in California, such as median income, house age, and number of rooms. The goal of the model is to predict the median house value (`MedHouseVal`) based on these features.

The code first loads the dataset using the `fetch_california_housing()` function from scikit-learn and converts it to a Pandas DataFrame using the `as_frame=True` option. It then splits the data into training and testing sets using the `train_test_split()` function, where 80% of the data is used for training and the remaining 20% is used for testing.

Next, the code trains a linear regression model using the training data and the `LinearRegression()` function from scikit-learn. It then evaluates the model's performance on the testing data using the mean squared error (MSE) metric, which measures the average squared difference between the predicted and actual target values.

This code can be helpful for predicting house prices using the California Housing dataset by providing a simple implementation of a linear regression model and a way to evaluate its performance. It can also serve as a starting point for more advanced modeling techniques and feature engineering to improve the model's accuracy.

how these code can be helpful in real-world for Predicting House Prices with California Housing Dataset

This code can be helpful in real-world applications for predicting house prices with the California Housing dataset by providing a simple implementation of a linear regression model and a way to evaluate its performance. The model can be used to make predictions on new, unseen data based on the housing features provided in the dataset, which can be helpful for real estate agents, homeowners, and investors.

For example, a real estate agent could use this model to estimate the price of a house in a particular district based on its features, such as the number of rooms, median income, and house age. Investors could also use this model to identify undervalued properties that have the potential to increase in value.

Additionally, the code can serve as a starting point for more advanced modeling techniques and feature engineering to improve the model's accuracy. Real-world applications may require more sophisticated modeling techniques and additional data sources, such as location-based data, to produce more accurate predictions.