

Hands-on Lab: Final Project: Generative AI for Data Science

Estimated Effort: 60 mins

Project Scenario

You have been employed as a Data Scientist by a consultancy firm. The firm has a client who is a used car dealer. They have a special feature on Ford cars and they want your firm to design a model that can predict the optimum quotation price for the cars in their lot. They provide you with sales data for the past few years. The dataset contains different features of the cars and the price they were sold at.

The tasks assigned to you are as follows.

1. There might be a few duplicate entries and a few missing values in the dataset. Data cleaning will be a part of the assignment.
2. You have to perform exploratory data analysis to draw keen insights on the data and determine the effect of different features on the price. Some specific requests by the client include:
 - a. Identify number of sales for each fuel type
 - b. Identify which transmission type has more price outliers
3. Compare the models with linear, polynomial and ridge regressions on single and multiple variables to find the best performing model
4. Perform a Grid Search on the Ridge regression model to identify the optimum hyperparameter for the model for best performance.

You decide to use Generative AI to create python codes that can help you analyse the data, determine the best features and create the prediction model as per requirement.

Disclaimer: This is a fictitious scenario created for the purpose of this project. The dataset being used is publicly available.

About the Dataset

This dataset contains used car sale prices for Ford cars. This is a public dataset available on the [Kaggle](#) website as [Ford Car Pricing Dataset](#) under the [CC0: Public Domain](#) license. The dataset has been slightly modified for the purpose of this project.

Attributes of this dataset have been explained below.

Variable	Description
model	Car model name
year	Year of car make
transmission	Type of transmission (Automatic, Manual or Semi-Auto)
mileage	Number of miles traveled
fuelType	The type of fuel the car uses (Petrol, Diesel, Hybrid, Electric, Other)
tax	Annual Tax payable in USD
mpg	Miles per Gallon that the car runs at
engineSize	Engine Size of the car
price	Price of car in USD

Code execution environment

To test the prompt-generated code, keep the Jupyter Notebook (in the link below) open in a separate tab in your web browser. The notebook has some setup instructions that you should complete now.

Jupyter-Lite Test Environment

Please note the lab environment above will only work on Windows (Google Chrome or Firefox browser). If you don't have a Windows system with either of these browsers, use the lab environment provided in the next lesson of the module.

The data set for this lab is available in the following URL.

URL = https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMSkillsNetwork-AI0271EN-SkillsNetwork/labs/v1/m3/data/used_car_price_analysis.csv

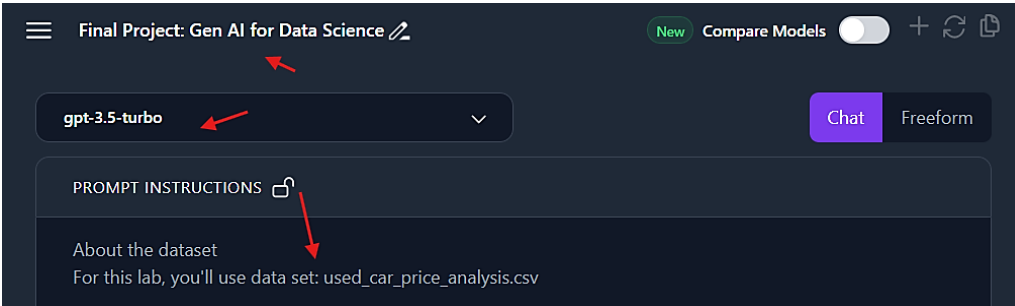
```
1 URL = "https://cf-courses-data.s3.us.cloud-object-storage.appdomain.cloud/IBMSkillsNetwork-AI0271EN-SkillsNetwork/labs/v1/m3/data/used_car_price_analysis."
```

Complete the setup in the Jupyter Notebook and then proceed further.

Important Note: All prompts that are made available have been hidden and the users are encouraged to first try to write their own prompts to create the solutions. Also, the prompts given as solutions have also been maintained as ones which will create generic code structures which you can modify according to the question at hand.

Importing the Dataset

You can begin by using the Generative AI model to create a python script that can load the dataset to a pandas dataframe. The dataset file already has the headers in the first row.



Write the prompt to generate the said code and test it in the JupyterLite environment. For verification of appropriate loading, include a step for printing the first 5 values of the loaded dataframe. ▼

Click here for the prompt

- Write a Python code that can perform the following tasks.
1. Read the CSV file, located on a given file path, into a pandas data frame, assuming that the first row of the file can be used as the headers for the data.
 2. Print the first 5 rows of the dataframe to verify correct loading.

```
1 Write a Python code that can perform the following tasks.
2 1. Read the CSV file, located on a given file path, into a pandas data frame, assuming that the first row of the file can be used as the headers for the data.
3 2. Print the first 5 rows of the dataframe to verify correct loading.
```

Write a Python code that can perform the following tasks.

1. Read the CSV file, located on a given file path, into a pandas data frame, assuming that the first row of the file can be used as the headers for the data.
2. Print the first 5 rows of the dataframe to verify correct loading.



This Python code reads a CSV file into a pandas data frame and then prints the first 5 rows of the dataframe to verify correct loading. Make sure to replace 'path_to_your_file/used_car_price_analysis.csv' with the actual file path where the CSV file is located.

```
1 import pandas as pd
2
3 # Specify the file path
4 file_path = 'path/to/used_car_price_analysis.csv'
5
6 # Read the CSV file into a pandas data frame
7 df = pd.read_csv(file_path)
8
9 # Print the first 5 rows of the dataframe
10 print(df.head())
```

Data Preparation

Data Cleaning

At this stage, it is required to clean up the data. As has been informed to you, the data may have missing values and duplicate entries. Write a prompt that performs the following tasks

1. Identifies the columns with missing values and fills the blank cells with average value of the columns.
2. Identifies and drops the duplicate entries from the data.

▼ Click here for prompt

Write a python code that performs the following tasks:

1. Identify the columns of a data frame with missing values.
2. Replace the missing values thus identified with mean value of the column.
3. Deletes the duplicate entries from the data frame.

1	Write a python code that performs the following tasks:
2	1. Identify the columns of a data frame with missing values.
3	2. Replace the missing values thus identified with mean value of the column.
4	3. Deletes the duplicate entries from the data frame.



This Python code defines a function that identifies columns with missing values in a dataframe, replaces the missing values with the mean of the column, and then deletes duplicate entries from the dataframe.

Write a python code that performs the following tasks:

1. Identify the columns of a data frame with missing values.
2. Replace the missing values thus identified with mean value of the column.
3. Deletes the duplicate entries from the data frame.

```
1 import pandas as pd
2
3 # Assume df is the pandas data frame
4
5 # Identify columns with missing values
6 columns_with_missing_values = df.columns[df.isnull().any()]
7
8 # Replace missing values with mean of the column using explicit assignment
9 for column in columns_with_missing_values:
10     df.fillna({column: df[column].mean()}, inplace=True)
11
12 # Delete duplicate entries
13 df.drop_duplicates(inplace=True)
```

Data Augmentation (optional)

Once cleaned, you may choose to augment this dataset with additional samples, created synthetically using [Mostly.ai](#).

python

```
1 df.to_csv('cleaned_data.csv', index=False)
```

Data Insights and Visualization

Write **prompts** that generate codes to perform the following actions.

1. Identify the **5 attributes** that have the **highest correlation** with the **price parameter**.

▼ Click here for the prompt

1

Write a python code that identifies the top 5 attributes with highest correlation with the target attribute in a data frame.



Write a python code that identifies the top 5 attributes with highest correlation with the target attribute in a data frame.

This Python code defines a function that calculates the correlation of all attributes with the target attribute in a dataframe, sorts them in descending order, and returns the top 5 attributes with the highest correlation.

```
1 import pandas as pd
2
3 # Assume df is the pandas data frame
4
5 # Select only numeric columns
6 numeric_df = df.select_dtypes(include='number')
7
8 # Calculate correlations with the target attribute
9 target_attribute = 'target'
10 correlations = numeric_df.corr()[target_attribute].abs().sort_values(ascending=False)
11 top_5_attributes = correlations[1:6]
12 print(top_5_attributes)
```

2. **Count** the number of cars under each **unique** value of **fuelType attribute**.

▼ Click here for the prompt

1

Write a python code that counts the number of entries in a dataframe with each unique value of a specific attribute.

```

1 import pandas as pd
2
3 # Assume df is the pandas data frame
4
5 # Specify the attribute for which unique value counts are needed
6 specific_attribute = 'attribute_name'
7
8 # Count the number of entries for each unique value of the specific attribute
9 value_counts = df[specific_attribute].value_counts()
10 print(value_counts)

```

3. Create a **Box plot** to determine whether cars with **automatic**, **manual** or **semi-auto** type of **transmission** have more **price outliers**. Use the **Seaborn** library for creating the plot.

▼ Click here for the prompt

1 Write a python code that creates a Seaborn box plot between a source and target attributes of a dataframe.

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Assume df is the pandas data frame
5 sns.boxplot(x='source_attribute', y='target_attribute', data=df)
6 plt.show()

```

4. Generate the **regression plot** between **mpg** parameter and the **price** to determine the **correlation type** between the two.

▼ Click here for the prompt

1 Write a python code that creates a Seaborn box plot between a source and target attributes of a dataframe.

```

1 import seaborn as sns
2 import matplotlib.pyplot as plt
3
4 # Assume df is the pandas data frame
5
6 # Create a regression plot between a source attribute and a target attribute
7 sns.regplot(x='source_attribute', y='target_attribute', data=df)
8 plt.show()

```

Model Development and Evaluation

Write prompts that generate codes to perform the following actions.

1. **Fit a linear regression model** to predict the **price** using the **feature mpg**. Then calculate the **R²** and **MSE values** for the model.

▼ Click here for the prompt

```
1 Write a python code to fit a linear regression model on a target variable on basis of a source variable. Calculate the R^2 values and MSE values for this model.
```

Write a python code to fit a linear regression model on a target variable on basis of a source variable. Calculate the R² values and MSE values for this model.

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score, mean_squared_error
3 import pandas as pd
4
5 # Assume df is the pandas data frame
6
7 # Define source and target variables
8 X = df[['source_variable']]
9 y = df['target_variable']
10
11 # Fit a linear regression model
12 model = LinearRegression()
13 model.fit(X, y)
14
15 # Make predictions
16 y_pred = model.predict(X)
17
18 # Calculate R^2 and MSE
19 r2 = r2_score(y, y_pred)
20 mse = mean_squared_error(y, y_pred)
21
22 print("R^2 Score:", r2)
23 print("Mean Squared Error:", mse)
```

2. Fit a **linear regression model** to predict the price using the following **set of features**.
year, mileage, tax, mpg and **engineSize**.
Calculate the **R²** and **MSE** values for this **model**.

▼ Click here for the prompt

```
1 Write a python code to fit a linear regression model on a target variable on basis of a set of source variables. Calculate the R^2 values and MSE values for this model.
```

Write a python code to fit a linear regression model on a target variable on basis of a set of source variables. Calculate the R² values and MSE values for this model.

```
1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import r2_score, mean_squared_error
3 import pandas as pd
4
5 # Assume df is the pandas data frame
6
7 # Define source and target variables
8 X = df[['source_variable_1', 'source_variable_2', 'source_variable_3']] # Add all relevant source variables
9 y = df['target_variable']
10
11 # Fit a linear regression model
12 model = LinearRegression()
13 model.fit(X, y)
14
15 # Make predictions
16 y_pred = model.predict(X)
17
18 # Calculate R^2 and MSE
19 r2 = r2_score(y, y_pred)
20 mse = mean_squared_error(y, y_pred)
21
22 print("R^2 Score:", r2)
23 print("Mean Squared Error:", mse)
```


3. For the **same set of features** as in the question above, create a **pipeline model** object that uses **standard scalar**, **second degree polynomial features** and a **linear regression model**.

Calculate the **R²** value and the **MSE value** for this **model**.

▼ Click here for the prompt

Write a **python code** that can perform the following tasks.

```
1 Write a python code that can perform the following tasks.
2
3 1. Create a pipeline with standard scalar, second degree polynomial features and linear regression model.
4 2. Fit this pipeline for a target variable using a set of source attributes from a dataframe.
5 3. Evaluate the R^2 and MSE values for the trained model.
```

Write a python code that can perform the following tasks.

1. Create a pipeline with standard scalar, second degree polynomial features and linear regression model.
2. Fit this pipeline for a target variable using a set of source attributes from a dataframe.
3. Evaluate the R² and MSE values for the trained model.

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.preprocessing import StandardScaler, PolynomialFeatures
3 from sklearn.linear_model import LinearRegression
4 from sklearn.metrics import r2_score, mean_squared_error
5 import pandas as pd
6
7 # Assume df is the pandas data frame
8
9 # Define source and target variables
10 X = df[['source_variable_1', 'source_variable_2', 'source_variable_3']] # Add all relevant source variables
11 y = df['target_variable']
12
13 # Create a pipeline with standard scalar, polynomial features, and linear regression
14 pipeline = Pipeline([
15     ('scaler', StandardScaler()),
16     ('poly_features', PolynomialFeatures(degree=2)),
17     ('linear_reg', LinearRegression())
18 ])
19
20 # Fit the pipeline
21 pipeline.fit(X, y)
22
23 # Make predictions
24 y_pred = pipeline.predict(X)
25
26 # Calculate R^2 and MSE
27 r2 = r2_score(y, y_pred)
28 mse = mean_squared_error(y, y_pred)
29
30 print("R^2 Score:", r2)
31 print("Mean Squared Error:", mse)
```

- For the same set of features, split the data into **training** and **testing** data parts. Assume testing part to be **20%**. Create and **fit** a **Ridge regression** object using the training data, set the **regularization parameter** to **0.1**, and calculate the **R²** using the **test** data.

▼ Click here for the prompt

```
1 Write a python code that can perform the following tasks.
2
3 1. Assuming that a subset of the attributes of a data frame are source attributes and one of the attributes is a target attribute, split the data into training and testing data assuming the testing data to be 20%.
4 2. Create and fit a Ridge regression model using the training data, setting the regularization parameter to 0.1.
5 3. Calculate the MSE and R^2 values for the Ridge Regression model using the testing data.
```

Write a python code that can perform the following tasks.

- Assuming that a subset of the attributes of a data frame are source attributes and one of the attributes is a target attribute, split the data into training and testing data assuming the testing data to be 20%.
- Create and fit a Ridge regression model using the training data, setting the regularization parameter to 0.1.
- Calculate the MSE and R² values for the Ridge Regression model using the testing data.

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.linear_model import Ridge
3 from sklearn.metrics import r2_score, mean_squared_error
4 import pandas as pd
5
6 # Assume df is the pandas data frame
7
8 # Define source and target variables
9 X = df[['source_variable_1', 'source_variable_2', 'source_variable_3']] # Add all relevant source variables
10 y = df['target_variable']
11
12 # Split the data into training and testing sets
13 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
14
15 # Create and fit a Ridge regression model
16 ridge_model = Ridge(alpha=0.1)
17 ridge_model.fit(X_train, y_train)
18
19 # Make predictions on the testing data
20 y_pred = ridge_model.predict(X_test)
21
22 # Calculate R^2 and MSE for the Ridge Regression model
23 r2 = r2_score(y_test, y_pred)
24 mse = mean_squared_error(y_test, y_pred)
25
26 print("R^2 Score:", r2)
27 print("Mean Squared Error:", mse)
```

5. Perform a **second order polynomial transform** on both the **training** data and **testing** data created for the question above. Create and **fit a Ridge regression object** using the modified training data, set the **regularisation** parameter to **0.1**, and calculate the **R²** and **MSE** utilising the **modified test** data.

▼ Click here for the prompt

```
1  Write a python code that can perform the following tasks.
```

- ```
2
3 1. Assuming that a subset of the attributes of a data frame are source attributes and one of the attributes is a target attribute, split the data into training and testing data assuming the testing data to be 20%.
4 2. Apply second degree polynomial scaling to the training and testing data.
5 3. Create and fit a Ridge regression model using the training data, setting the regularization parameter to 0.1.
6 4. Calculate the MSE and R^2 values for the Ridge Regression model using the testing data.
```

Write a python code that can perform the following tasks.

1. Assuming that a subset of the attributes of a data frame are source attributes and one of the attributes is a target attribute, split the data into training and testing data assuming the testing data to be 20%.
2. Apply second degree polynomial scaling to the training and testing data.
3. Create and fit a Ridge regression model using the training data, setting the regularization parameter to 0.1.
4. Calculate the MSE and R<sup>2</sup> values for the Ridge Regression model using the testing data.

```
1 from sklearn.model_selection import train_test_split
2 from sklearn.preprocessing import PolynomialFeatures
3 from sklearn.linear_model import Ridge
4 from sklearn.metrics import r2_score, mean_squared_error
5 import pandas as pd
6
7 # Assume df is the pandas data frame
8
9 # Define source and target variables
10 x = df[['source_variable_1', 'source_variable_2', 'source_variable_3']] # Add all relevant source variables
11 y = df['target_variable']
12
13 # Split the data into training and testing sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
15
16 # Apply second degree polynomial scaling to the training and testing data
17 poly = PolynomialFeatures(degree=2)
18 X_train_poly = poly.fit_transform(X_train)
19 X_test_poly = poly.transform(X_test)
20
21 # Create and fit a Ridge regression model
22 ridge_model = Ridge(alpha=0.1)
23 ridge_model.fit(X_train_poly, y_train)
24
25 # Make predictions on the testing data
26 y_pred = ridge_model.predict(X_test_poly)
27
28 # Calculate R^2 and MSE for the Ridge Regression model
29 r2 = r2_score(y_test, y_pred)
30 mse = mean_squared_error(y_test, y_pred)
31
32 print("R^2 Score:", r2)
33 print("Mean Squared Error:", mse)
```

6. In the question above, perform a **Grid Search** on **ridge regression** for a set of values of **alpha** {0.01, 0.1, 1, 10, 100} with **4-fold cross validation** to find the optimum value of **alpha** to be used for the **prediction model**.

▼ Click here for the prompt

```
1 Write a python code that can perform the following tasks.
2
3 1. Assuming that a subset of the attributes of a data frame are source attributes and one of the attributes is a target attribute, split the data into training and testing data assuming the testing data to be 20%.
4 2. Apply second degree polynomial scaling to the training and testing data.
5 3. Create and fit a Grid search on Ridge regression with cross validation using the training data, for a set of values of the parameter alpha.
6 4. Calculate the MSE and R^2 values for the Ridge Regression model using the testing data.
```

Write a python code that can perform the following tasks.

1. Assuming that a subset of the attributes of a data frame are source attributes and one of the attributes is a target attribute, split the data into training and testing data assuming the testing data to be 20%.
2. Apply second degree polynomial scaling to the training and testing data.
3. Create and fit a Grid search on Ridge regression with cross validation using the training data, for a set of values of the parameter alpha.
4. Calculate the MSE and R^2 values for the Ridge Regression model using the testing data.

```
1 from sklearn.model_selection import train_test_split, GridSearchCV
2 from sklearn.preprocessing import PolynomialFeatures
3 from sklearn.linear_model import Ridge
4 from sklearn.metrics import r2_score, mean_squared_error
5 import pandas as pd
6
7 # Assume df is the pandas data frame
8
9 # Define source and target variables
10 X = df[['source_variable_1', 'source_variable_2', 'source_variable_3']] # Add all relevant source variables
11 y = df['target_variable']
12
13 # Split the data into training and testing sets
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
15
16 # Apply second degree polynomial scaling to the training and testing data
17 poly = PolynomialFeatures(degree=2)
18 X_train_poly = poly.fit_transform(X_train)
19 X_test_poly = poly.transform(X_test)
20
21 # Create Ridge regression model
22 ridge = Ridge()
23 parameters = {'alpha': [0.1, 1, 10]} # Define set of alpha values to search
24
25 # Perform Grid Search with Cross Validation
26 ridge_cv = GridSearchCV(ridge, parameters, cv=5)
27 ridge_cv.fit(X_train_poly, y_train)
28
29 # Make predictions on the testing data
30 y_pred = ridge_cv.predict(X_test_poly)
31
32 # Calculate R^2 and MSE for the Ridge Regression model
33 r2 = r2_score(y_test, y_pred)
34 mse = mean_squared_error(y_test, y_pred)
35
36 print("R^2 Score:", r2)
37 print("Mean Squared Error:", mse)
```

# Conclusion

Congratulations! You have completed this guided project on using Generative AI for different data science tasks.

By the end of this project, you are now capable of using Generative AI for the tasks of:

- **Data preparation:** cleaning, transforming and augmentation
- **Data analysis:** drawing insight, creating visualizations
- **Model development:** creating simple as well as complex prediction models
- **Model refinement:** found the optimum model using Grid Search

## Author(s)

[Abhishek Gagneja](#)