

## Homework 2 (HW2)

By the end of this homework, we expect you to be able to:

- Preprocess data and make it amenable to statistical analysis and machine learning models;
- Train and test out-of-the-box machine learning models in `sklearn`;
- Carry out simple logistic regression and decision tree analysis;
- Use propensity score matching to estimate treatment effects;

## Important Dates

- Homework release: Fri 15 Nov 2024
- Homework due: Fri 29 Nov 2024, 23:59
- Grade release: Mon 09 Dec 2024

## Some rules

1. We have provided 7000+ comments in the code cells that you need to fill out with your solutions. For some questions, we have also provided Your response comments, where you should provide a textual answer.
2. You are allowed to use any built-in Python library that is included in the requirements.txt for this homework. If you use any additional library, this may complicate the grading process, and we reserve the right to penalize your grade for unnecessary complexity of the solution. All the questions can be solved with the libraries in requirements.txt.
3. Please write all your comments in English, and use meaningful variable names in your code. Your repo should have a single notebook (plus the required data files) in the master/main branch. If there are multiple notebooks present, we will not grade anything.
4. We will not run your notebook for you! Rather, we will grade it as is, which means that only the results contained in your evaluated code cells will be considered, and we will not see the results in unevaluated code cells. Thus, be sure to hand in a fully-run and evaluated notebook. In order to check whether everything looks as intended, you can check the rendered notebook on the GitHub website once you have pushed your solution there.
5. In continuation to the previous point on additional library, interactive plots, such as those generated using `plotly`, should be strictly avoided!

### A Note on Using Language Models (LMs)

If you try hard enough, you will likely get away with cheating. Fortunately, our job is not to police, but rather to educate! So, please consider the following:

Presumably, you are taking this course to learn something! LMs are not always right (they often fall in silly ways). This course should prepare you to detect when they are wrong! Some of the TAs on this course literally published many works on detecting machine-generated text.

Here LLM includes but not limited to chatbots like ChatGPT, coding assistants like Copilot. Do not even use them to prettify your code or correct English. If you are caught using LLMs, you will be reported to the instructor and subject to the consequences.

## Grading

- The homework has a total of 100 points, distributed as follows:
  - Part 1: Data Preprocessing (20 points)
  - Part 2: Linear Regression (30 points)
  - Part 3: Supervised Learning (40 points)
  - Part 4: Propensity Score Matching (10 points)

## Context

Within EPFL's master program, you are excited to start an internship as a data scientist. After rounds of interviews, you have been selected to work with the biggest car dealership in Switzerland!

Your mentor at the company Tim, has explained to you that the company is interested in a pricing model for used cars.

- Tim: "We have a lot of used cars in our inventory, and we need to determine the price at which we should sell these cars. We have some ideas about the factors that influence the price of a used car, but so far we have just been using our experience and intuition to determine the price of a used car. Sometimes it works, but probably we can do better and a more data-driven approach would also help our new employees in the sales team as they have less experience."
- You: "That sounds like a great project! What kind of data do we have?"
- Tim: "We sell all kinds of cars here, but maybe we can start with a specific brand and model. For example, the Toyota Corolla is the best-selling car worldwide in 2023, and we have a lot of data on it. We can start by analyzing the data on used Toyota Corolla cars. If it works well, we can extend the analysis to other brands."

The dataset contains the following columns:

- **Age**: Age of the car in months.
- **Mileage**: Number of distance the car has been driven, (km or miles)
- **FuelType**: Fuel type of the car (Petrol, Diesel, or CNG)
- **HP**: Horsepower
- **MetColor**: Is the color of the car metallic? (Yes=1, No=0)
- **Automatic**: Is the car automatic? (Yes=1, No=0)
- **CC**: Cylinder volume in cubic centimeters
- **Doors**: Number of doors
- **Weight**: Weight of the car in kilograms
- **Price**: Price of the car in euros

## Data

The data is provided in the `data` folder and it contains the following 3 csv files:

- `Task1-2_ToyotaCorolla.csv` for Part 1 and Part 2
- `Task3_ToyotaCorolla_sales_3months.csv` for Part 3
- `Task4_ToyotaCorolla_discount_sales` for Part 4

You should not use any other data source for this homework.

For some questions, you might need to slightly modify the data. But overall, you should avoid making any major changes to the data, which may affect your analysis.

## References:

The data is based on the ToyotaCorolla dataset from the UCI Machine Learning Repository [here](#). We have made some modifications to the original dataset, so please use the data provided in the `data` folder in the course repo.

## Task 1 (20 pts) - Get to know the data

In [13]:

```
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import statsmodels.api as sm
import random
import sklearn

# fix random seed for reproducibility
np.random.seed(42)
random.seed(42)
```

**1.1 (2 pts):** Load the data from the file `Task1-2_ToyotaCorolla-raw.csv` into a pandas DataFrame. Display the first 5 rows of the DataFrame. Hint: A naive loading of the data will raise an error. You will need to figure out how to load the data correctly. (Hint: localise which row is causing the error)

In [14]:

```
# Load the data

# data_df = pd.read_csv('data/Task1-2_ToyotaCorolla.csv') this will fail because the file is not in
# The best way is to fix the line 33 in the file by either manually editing it or writing a function

try:
    data_df = pd.read_csv('data/Task1-2_ToyotaCorolla-raw.csv')
except Exception as e:
    print(e)

# If the students skip the bad line with the following or similar way, they will only get 1/2 point

data_df = pd.read_csv('data/Task1-2_ToyotaCorolla-raw.csv', on_bad_lines='skip')

# Display the first 5 rows of the data
print(data_df.head())
```

Error tokenizing data. C error: Expected 11 fields in line 33, saw 12

```
   Price   Age   Mileage   FuelType   HP   MetColor   Automatic   CC   Doors   \
0  13500.0  23.0  46986.0   Diesel   90.0         1.0         0.0  2000.0   3.0
1  11878.0  20.0  72937.0   Diesel   90.0         1.0         0.0  2000.0   3.0
2  12059.0  24.0  41711.0   Diesel   90.0         1.0         0.0  2000.0   3.0
3  12914.0  26.0  48000.0   Diesel   90.0         0.0         0.0  2000.0   3.0
4  11878.0  30.0  38500.0   Diesel   90.0         0.0         0.0  2000.0   3.0
```

```
   Weight   Currency
0  1165.0      EURO
1  1165.0      CHF
2  1165.0      CHF
3  1165.0      CHF
4  1170.0      CHF
```

**1.2 (2 pts):** Check if there are nan values in the DataFrame. If there are, try to find out which row is problematic and fix it. If you can't fix it, drop the row.

In [15]:

```
# get nan-index
nan_index = data_df.isnull().any(axis=1)
print(data_df[nan_index])

# It seems that on that row the delimiter is not correct, it uses a semicolon instead of a comma, so
# fix it by replacing the semicolon with a comma manually or by writing a function to do it

# drop the nan row after we manually fixed it
data_df = data_df.dropna()
```

```
   Price   Age   Mileage   FuelType   HP   MetColor   Automatic   CC   Doors   \
905  8423.0;68;58860.0;Petrol;110;1;0;1600;3;1055;CHF   NaN   NaN   NaN
906  HP   MetColor   Automatic   CC   Doors   Weight   Currency
      NaN   NaN   NaN   NaN   NaN   NaN   NaN
```

In [16]:

```
nan_index = data_df.isnull().any(axis=1)
print(data_df[nan_index])

Empty DataFrame
Columns: [Price, Age, Mileage, FuelType, HP, MetColor, Automatic, CC, Doors, Weight, Currency]
Index: []
```

**1.3 (4 pts): Compute the mean, median of the `Price` column.**

- Compute the mean and median of the `Price` column. If you encounter error, try to understand why this error is happening and propose a solution.

Hint: Is all values in the `Price` column numerical?

In [123]:

```
# directly apply mean will not work because the column is not numeric due to the ; symbol
try:
    data_df['Price'].mean()
except Exception as e:
    print(e)

# List the unique values in the column to see what is wrong
# print(data_df['Price'].unique())

# 2 issues:
# 1. there is some values with ; symbol (-1 point if not found this issue)
# 2. there use comma as decimal separator instead of dot (-1 point if not found this issue)
```

In [118]:

```
def convert_raw_price_to_price(raw_price:str)->float:
    if raw_price.endswith(';'):
        raw_price = raw_price.replace(';','')
    # some numbers have comma as thousand separator
    raw_price = raw_price.replace(',','')
    return float(raw_price)

# Convert the price column to float
data_df['Price'] = data_df['Price'].apply(convert_raw_price_to_price)

mean = data_df['Price'].mean()
median = data_df['Price'].median()

print(f'Mean: {mean}, Median: {median}')
```

Mean: 10738662038331.08, Median: 8595.0

In [119]:

```
# The huge difference implies that there may be an outlier that is skewing the mean
# Let's sort the prices and see if there are any outliers
sorted_prices = data_df['Price'].sort_values()
print(sorted_prices)
```

```
100      3.580000e+03
1047     3.801000e+03
392      3.844000e+03
1368     4.090000e+03
191      4.183000e+03
      ...
178      2.250000e+04
109      2.677900e+04
110      2.701600e+04
108      2.807400e+04
107      3.539500e+04
Name: Price, Length: 1435, dtype: float64
```

In [120]:

```
# remove the last value as it is an outlier
data_df = data_df[data_df['Price'] != sorted_prices.iloc[-1]]
new_mean = data_df['Price'].mean()
new_median = data_df['Price'].median()
print(f'New Mean: {new_mean}, New Median: {new_median}') # If the student didn't compute the new me

New Mean: 9431.738438563458, New Median: 8595.0
```

**1.4 (4 pts): Convert Units**

You notice that some prices are in CHF (Swiss Francs), while others are in EUR (Euros) or GBP (British Pounds). Additionally, for cars priced in GBP, the mileage is in miles rather than kilometers.

For consistency, convert all prices to CHF and all distances to kilometers.

- Exchange rates:
  - 1 CHF = 1.05 EUR
  - 1 GBP = 1.15 CHF
  - 1 mile = 1.61 km

Make the following conversions:

1. Convert prices in EUR or GBP to CHF, rounding to the nearest integer.
2. Convert distances in miles (for GBP cars) to kilometers, rounding to the nearest integer.
3. Drop the 'Currency' column.
4. Calculate the min, mean, median and max of the 'Price' and 'Distance' columns after the conversion.

In [124]:

```
# Load the cleaned data

data_df = pd.read_csv('data/Task1-2_ToyotaCorolla-clean.csv')

# shuffle the data and sample it
data_df = data_df.sample(frac=1)

data_df.to_csv('data/Task1-2_ToyotaCorolla-clean-shuffled.csv', index=False)

# Convert the prices to CHF, round to the nearest integer
data_df['Price'] = np.where(data_df['Currency'] == 'EUR', data_df['Price']*1.05, data_df['Price'])
data_df['Price'] = np.where(data_df['Currency'] == 'GBP', data_df['Price']*1.15, data_df['Price'])
data_df['Price'] = np.round(data_df['Price'])

# Convert the distance to kilometers
data_df['Mileage'] = np.where(data_df['Currency'] == 'GBP', data_df['Mileage']*1.61, data_df['Mileage'])
data_df['Mileage'] = np.round(data_df['Mileage'])

# Drop the 'Currency' column
data_df = data_df.drop(columns=['Currency'])

# Calculate the min, mean, median and max of the 'Price' and 'Distance' columns after the conversion
print(data_df[['Price', 'Mileage']].agg(['min', 'mean', 'median', 'max']))

# the correct values should be below, if Mileage/Price column does not have the correct values, the
```

```
   Price   Mileage
min  3758.000000  1.000000
mean  9423.536212  68533.259749
median  8595.000000  63309.500000
max  28074.000000  243000.000000
```

**1.5 (2 pts): Analyze Average Price**

A. Print the average price for each fuel type. Determine which fuel type has the highest average price.

B. Print the average price for different numbers of doors. Determine which number of doors has the highest average price.

In [22]:

```
print(data_df.groupby('FuelType')['Price'].mean())
print(data_df.groupby('Doors')['Price'].mean())

# 1 point should be deducted if the student does not use groupby function, even though the result is
```

```
FuelType      8383.235294
CNG           9878.613355
Diesel        9881.722310
Name: Price, dtype: float64
Doors
2      6997.500000
3      8898.673633
4      8604.384058
5      10082.023442
Name: Price, dtype: float64
```

**1.6 (2 pts): Relationship Between Car Age and Price**

It is intuitive that an older car tends to be cheaper, and a car with more mileage might also be less expensive.

To explore this intuition, create two scatter plots:

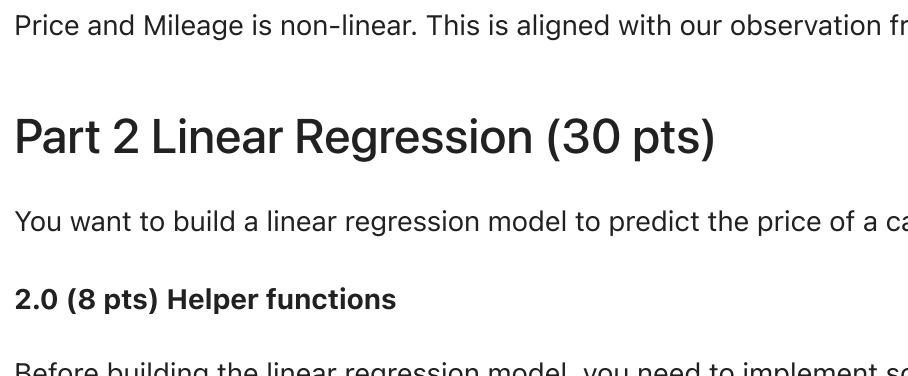
1. Car Age vs Price
2. Mileage vs Price

In [23]:

```
plt.figure(figsize=(5,5))
plt.scatter(data_df['Age'], data_df['Price'])
plt.title('Price vs Age')
plt.xlabel('Age(years)')
plt.ylabel('Price(CHF)')
plt.show()
```

```
plt.figure(figsize=(5,5))
plt.scatter(data_df['Mileage'], data_df['Price'])
plt.title('Price vs Mileage')
plt.xlabel('Mileage(km)')
plt.ylabel('Price(CHF)')
plt.show()
```

# we don't grade based on the style of the image but if there is missing labels, title, or not unit



**1.7 (4 pts): Correlation Between Price and Mileage**

The relationship between car price and mileage appears non-linear, with a steeper price drop initially followed by a flatter curve.

A. (2 pts) Calculate both the Pearson and Spearman correlations between the price of the car and the distance driven.

In [24]:

```
# Calculate Pearson and Spearman correlations between Price and Mileage
pearson_corr = data_df[['Price', 'Mileage']].corr(method='pearson')['Price']['Mileage']
spearman_corr = data_df[['Price', 'Mileage']].corr(method='spearman')['Price']['Mileage']

print("Pearson Correlation:\n", pearson_corr)
print("Spearman Correlation:\n", spearman_corr)

# if the correlation values are not correct, 1 point should be deducted

Pearson Correlation:
-0.563456654608847
Spearman Correlation:
-0.605359801804244
```

B. (2 pts) Which correlation value is higher? Does this result align with your expectations?

Your Response: Spearman correlation is higher than Pearson correlation, which means that the relationship between Price and Mileage is non-linear. This is aligned with our observation from the scatter plot.

## Part 2 Linear Regression (30 pts)

You want to build a linear regression model to predict the price of a car based on the features you have.

### 2.0 (8 pts) Helper functions

Before building the linear regression model, you need to implement some helper functions.

Implement the `accuracy`, `precision`, `recall` and `f1_score` functions.

1. These functions should take in the true labels (`np.array`) and the predicted labels (`np.array`) and return the corresponding metric.
2. They should follow the convention that the positive class is 1 and the negative class is 0.
3. Apply the functions to the following data:

```
true_labels = np.array([1, 0, 1, 1, 0, 1, 0, 0, 1, 0])
predicted_labels = np.array([1, 1, 1, 1, 0, 0, 1, 0, 1, 0])

# Compare the results with the implementation in sklearn and see if they match.
```

In [25]:

```
# Implement the accuracy, precision, recall, and f1_score functions using vectorized operations with
# sklearn.metrics

def accuracy_fn(y_true, y_pred):
    return np.mean(y_true == y_pred)

def precision_fn(y_true, y_pred):
    tp = np.sum(np.logical_and(y_true == 1, y_pred == 1))
    fp = np.sum(np.logical_and(y_true == 0, y_pred == 1))
    if tp + fp == 0:
        return 0.0 # Define precision as 0 when no positive predictions
    return tp / (tp + fp)

def recall_fn(y_true, y_pred):
    tp = np.sum(np.logical_and(y_true == 1, y_pred == 1))
    fn = np.sum(np.logical_and(y_true == 1, y_pred == 0))
    if tp + fn == 0:
        return 0.0 # Define recall as 0 when no positive true labels
    return tp / (tp + fn)

def f1_score_fn(y_true, y_pred):
    prec = precision_fn(y_true, y_pred)
    rec = recall_fn(y_true, y_pred)
    if prec + rec == 0:
        return 0.0 # Define F1-score as 0 when both precision and recall are 0
    return 2 * prec * rec / (prec + rec)
```

```
# Comments:
# If not using vectorized operations, 4 pts will be deducted
# If not using safe division, 2 pts will be deducted

# Test the functions
true_labels = np.array([1, 0, 1, 1, 0, 1, 0, 0, 1, 0])
predicted_labels = np.array([1, 1, 1, 1, 0, 0, 1, 0, 1, 0])

print("Accuracy:", accuracy_fn(true_labels, predicted_labels))
print("Precision:", precision_fn(true_labels, predicted_labels))
print("Recall:", recall_fn(true_labels, predicted_labels))
print("F1-Score:", f1_score_fn(true_labels, predicted_labels))
```

# see if the results are correct with sklearn

```
import sklearn.metrics

print("Accuracy:", sklearn.metrics.accuracy_score(true_labels, predicted_labels))
print("Precision:", sklearn.metrics.precision_score(true_labels, predicted_labels))
print("Recall:", sklearn.metrics.recall_score(true_labels, predicted_labels))
print("F1-Score:", sklearn.metrics.f1_score(true_labels, predicted_labels))
```

# 4 points should be removed if the student does not consider the case where the denominator is 0

```
Accuracy: 0.7
Precision: 0.6666666666666666
Recall: 0.8
F1-Score: 0.7272727272727272
Accuracy: 0.7
Precision: 0.6666666666666666
Recall: 0.8
F1-Score: 0.7272727272727272
```

**2.1 (6 pts) Preprocess the Data**

To prepare your data for building a linear regression model, complete the following steps:

A. (1 pts) Convert the categorical variables to one-hot encoding using the `pd.get_dummies()` function, how many columns do you have after the one-hot encoding? (P.S. You may want to avoid introducing multicollinearity with one-hot encoding, what should you do to avoid this?)

In [41]:

```
data_df = data_df.sort_index()
```

In [42]:

```
data_df_processed = pd.get_dummies(data_df, columns=['FuelType', 'MetColor', 'Automatic', 'Doors'],
                                drop_first=True)

# If they considered the 'Doors' column not as categorical, it's also correct.
data_df_processed_door_non_categorical = pd.get_dummies(data_df, drop_first=True)

# If drop_first not set, they will lose 1 point

print(f'Columns in the processed data: {len(data_df_processed.columns)} or {len(data_df_processed_door_non_categorical.columns)}

Columns in the processed data: len=13 or len=11
```

In [43]:

```
print(data_df_processed.columns)
print(data_df_processed_door_non_categorical.columns)
```

```
Index(['Price', 'Age', 'Mileage', 'HP', 'CC', 'Weight', 'FuelType_Diesel',
      'FuelType_Petrol', 'MetColor_1', 'Automatic_1', 'Doors_3', 'Doors_4',
      'Doors_5'],
      dtype='object')
Index(['Price', 'Age', 'Mileage', 'HP', 'MetColor', 'Automatic', 'CC', 'Doors',
      'Weight', 'FuelType_Diesel', 'FuelType_Petrol', 'MetColor_1',
      'Automatic_1', 'Doors_3', 'Doors_4', 'Doors_5'],
      dtype='object')
```

B. (1 pts) Split the data into features (X) and target (y) variables. The target variable is the 'Price' column. Then split the data into train test sets using a 80-20 split. Use `random_state=42` for reproducibility. How many samples are in the training set and how many samples are in the test set?

In [44]:

```
X = data_df_processed.drop(columns=['Price'])
y = data_df_processed['Price']

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

print(f'No. of training samples: {len(X_train)} No. of testing samples: {len(X_test)}')
```

No. of training samples: 1148 No. of testing samples: 288

C. (1 pts) Why do we split the data into only train-test sets but not train-validation-test sets?

Your Response: We split the data into train-test sets to train the model on the training set and evaluate it on the test set. The validation set is used to tune the hyperparameters of the model, and we don't need it here because we are not tuning any hyperparameters.

D. (1 pts) Standardize the Features: Use `StandardScaler` from `sklearn.preprocessing` and then add a constant column using `sm.add_constant()`. Print the average and standard deviation of the training set after standardization.

In [67]:

```
from sklearn.preprocessing import StandardScaler
import statsmodels.api as sm

scaler_train = StandardScaler()
scaler_train.fit(X_train)

X_train_scaled = pd.DataFrame(scaler_train.transform(X_train), index=X_train.index, columns=X_train.columns)
X_test_scaled = pd.DataFrame(scaler_train.transform(X_test), index=X_test.index, columns=X_test.columns)

# the student must use the scaler from the training set to scale the test set
# if they use a new scaler for the test set, 1 point should be deducted

# add constant column to both train and test sets
X_train_scaled_const = sm.add_constant(X_train_scaled)

X_test_scaled_const = sm.add_constant(X_test_scaled)

# print the avg and std of the training set
print("Training set mean: ", np.mean(X_train_scaled_const, axis=0))
print("Training set std: ", np.std(X_train_scaled_const, axis=0))
```

```
Training set mean: const      1.000000e+00
Age          1.485455e-16
Mileage      -1.485455e-16
HP           3.868373e-17
CC           -2.963174e-16
Weight       -6.854757e-16
FuelType_Diesel  1.005777e-17
FuelType_Petrol  2.321824e-17
MetColor_1   -9.739674e-17
Automatic_1   4.620470e-18
Doors_3       12.82949   437.519   -0.341   0.734   -1007.431   799.441
Doors_4       -20.1629   262.086   -0.077   0.346   -480.085   85.655
Doors_5       -214.0475   441.735   -0.486   0.627   -1081.557   651.862
dtype: float64
Training set std: const      0.0
```

```
Age          1.0
Mileage      1.0
HP           1.0
CC           1.0
Weight       1.0
FuelType_Diesel  1.0
FuelType_Petrol  1.0
MetColor_1   1.0
Automatic_1   1.0
Doors_3       1.0
Doors_4       1.0
Doors_5       1.0
dtype: float64
```

E. (2 pts) Should we first standardize the data and then split it into train and test sets or vice versa? why?

Your Response: We should first split the data into train and test sets and then standardize the features. This is because we want to avoid data leakage from the test set to the training set. If we standardize the data before splitting, the mean and standard deviation of the entire dataset will be used to standardize the data, which is incorrect.

### 2.2 (10 pts) Train and Evaluate the Linear Regression Model

To train a linear regression model using and evaluate its performance, follow these steps:

1. (2 pts) Train a linear regression model on the training dataset using `sm.OLS` from `statsmodels`, print the summary of the model using `model.summary()`.

In [68]:

```
logistic_regression_model = sm.OLS(X_train_scaled_const).fit()

print(logistic_regression_model.summary())

# normally the student should've got R-squared = 0.850 and Adj. R-squared = 0.848, but it's also possible
# as long as the values are between 0.82-0.88, it's acceptable
```

```
=====
                        OLS Regression Results
=====
Dep. Variable:          Price   R-squared:          0.850
Model:                  Least Squares   F-statistic:    536.0
Date:                   Tue, 83 Dec 2024   Prob (F-statistic): 0.00
Time:                   22:37:50   Log-Likelihood:   -9083.5
No. Observations:       1148   AIC:              1.963e+04
DF Residuals:           1145   BIC:              1.970e+04
Covariance Type:        nonrobust
=====
```

```
coef      std err      t      P>|t|      [0.025      0.975]
-----
const          9441.4983      36.725      257.086      0.000      9369.442      9513.555
Age            -0.97119732      51.805      -0.38605      0.000      -2073.618      -1870.328
Mileage         2.677900e+04      52.696      0.50846      0.000      -606.951      -600.165
HP              476.6196      94.331      5.053      0.000      291.538      661.701
CC             -477.0432      111.211      -4.290      0.000      -695.246      -258.841
Weight         -6.854757e-16      75.770      0.000      1.000      1084.292      1381.621
FuelType_Diesel  1.005777e-17      170.847      2.683      0.009      189.496      779.920
FuelType_Petrol  2.321824e-17      106.631      3.593      0.000      164.325      582.758
MetColor_1      -12.82949      37.119      0.346      0.730      -60.085      85.655
Automatic_1      4.620470e-18      39.024      0.242      0.809      -86.010      90.125
Doors_3          12.82949      437.519      -0.341      0.734      -1007.431      799.441
Doors_4          -20.1629      262.086      -0.077      0.939      -534.391      494.065
Doors_5          -214.0475      441.735      -0.486      0.627      -1081.557      651.862
=====
```

```
Omnibus:          171.243      Durbin-Watson:      2.076
Prob(Omnibus):    0.000      Jarque-Bera (JB):      2376.654
Skew:              -0.046      Prob(JB):              0.000
Kurtosis:          10.039      Cond. No.              34.4
=====
```

Notes: [1] Standard Errors assume that the covariance matrix of the errors is correctly specified.

2. (2 pts) Evaluate the model on the test dataset using the square root of the mean squared error (RMSE) metric. A. Report the RMSE value.

B. Your BOSS wants to know how far off the model's predictions are from the actual price of the car. What would you tell him? Give a number and explain how you got it.



## 2.5 (2 pts): Adding an Inverse Mileage Term

From the previous scatter plot, the relationship between car price and mileage appears non-linear, with a steep price drop initially and then a flattening. A suitable approach to model this behavior is by incorporating an inverse term of mileage.

- Add the inverse mileage term to the model and retrain it using the code provided. Print the model summary and interpret the effect of the inverse mileage term.

```
In [ ]: X['1/Mileage_inverse'] = 1/X['Mileage']

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)

scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
X_train = sm.add_constant(X_train)
X_test = sm.add_constant(X_test)

logistic_regression_model = sm.OLS(y_train, X_train).fit()
print(logistic_regression_model.summary())

y_pred = logistic_regression_model.predict(X_test)
r2 = r2_score(y_test, y_pred)

print("R2 score on test dataset: ", r2)

# comment:
# The student can report either the R2 score showed in the summary or the one calculated using the
# The R2 score should be pretty much the same as the one without the Inverse feature. This means the
# is not bringing any additional information to the model.

=====
OLS Regression Results
=====
Dep. Variable:      Price      R-squared:      0.850
Model:              OLS      Adj. R-squared:    0.849
Method:             Least Squares      F-statistic:    496.1
Date:               Tue, 03 Dec 2024      Prob (F-statistic): 0.000
Time:               22:41:48      Log-Likelihood: -9801.7
No. Observations:   1148      AIC:              1.963e+04
DF Residuals:       1134      BIC:              1.970e+04
DF Model:           13
Covariance Type:    nonrobust
=====
              coef      std err      t      P>|t      [0.025      0.975]
-----
const          9441.4983    36.685    257.367    0.000    9369.520    9513.476
x1             -1947.6675    51.769    -37.641    0.000   -2076.241   -1873.094
x2             -511.5765    52.819    -9.687    0.000   -615.310   -408.043
x3             472.3240    94.256     5.011    0.000    287.389    657.259
x4             -477.8449    111.091    -4.301    0.000   -695.811   -259.878
x5             1166.2577    76.053    15.343    0.000   1017.636    1316.079
x6             441.3814    170.670     2.586    0.010    106.516    776.247
x7             373.8717    106.515     3.510    0.000    164.083    582.860
x8             13.5550    37.081     0.366    0.715   -59.200    86.310
x9             -12.8931    30.406    -0.330    0.741   -85.464    63.677
x10            -150.5418    437.042    -0.344    0.731   -1008.044    706.961
x11            -17.8593    261.884    -0.068    0.946   -531.533    495.815
x12            -219.9402    441.262    -0.498    0.618   -1085.722    645.862
x13             -70.1493    37.611    -1.865    0.062   -143.943    3.645
=====
Durrbin-Watson:          167.538
Jarque-Bera (JB):        0.000
Prob(JB):                 0.000
Skewness:                 -0.001
Kurtosis:                 9.837
Cond. No.:                34.4
=====
```

Notes:  
[1] Standard Errors assume that the covariance matrix of the errors is correctly specified.  
R2 score on test dataset: 0.8132778753638819

## Part 3 Supervised Learning (40 pts)

After completing your analysis, you're satisfied with the results. You handed the Jupyter notebook over to your mentor.

(Fun fact: The name "Jupyter" is derived from Julia, Python, and R—three programming languages that the platform was originally designed for.)

Your mentor Tim is very impressed with your work and asks you the following question:

"This looks great! It will be very useful for our sales team. While looking at the results, I realized that there might be one thing that we can improve. For companies like us, it is important to sell the cars quickly. If we are patient, we might be able to sell the car for a higher price, but that's not always the best strategy. We need to consider the maintenance costs for the car, the cash flow and the fact that the price of the car decreases over time."

He then continues: "Three months is a sweet spot for us. If we can sell the car within the first three months, it is great. If not, it is worth considering lowering the price to sell it faster and increase our cash flow. I can ask Ivan from Sales to collect data in the last few months on whether the car was sold within the first three months or not. This would be great if you could have a model that tells us if the car will be sold in the first three months or not."

This sparks your interest, and soon Ivan has provided you with the new data containing an additional column `sold_within_3_months` which is a binary variable indicating whether the car was sold within the first three months or not.

Note: The data for this part is in the file `Task3.ToyotaCorolla_sales_3months.csv` and it has already unified the currency and distance units.

```
In [84]: data_df = pd.read_csv('data/Task3.ToyotaCorolla_sales_3months.csv', index_col=0)
print(data_df.head())

   Price  Age   KM FuelType  HP  MetColor  Automatic  CC  Doors  Weight  \
0  13500   23  46986   Diesel   90         1         0   2000      3   1165
1  13750   23  72937   Diesel   90         1         0   2000      3   1165
2  13950   24  41711   Diesel   90         1         0   2000      3   1165
3  14950   26  48000   Diesel   90         0         0   2000      3   1165
4  13750   30  38500   Diesel   90         0         0   2000      3   1165

   sold_within_3_months
0                        0
1                        0
2                        0
3                        0
4                        0
```

### 3.1 (2 pts): Preprocess the Data

- (1 pts) How many cars in the dataset were sold in the first three months, and how many were not?

```
In [85]: print(data_df['sold_within_3_months'].value_counts())

sold_within_3_months
0      880
1       56
Name: count, dtype: int64
```

- (1 pts) Preprocess the categorical variables to one-hot encoding using the `pd.get_dummies()` function.

```
In [86]: # 7000
data_df = pd.get_dummies(data_df, columns=['FuelType', 'MetColor', 'Automatic', 'Doors'], drop_first=True)
# comment:
# deduct 1 point if not using drop_first=True
```

### 3.2 (20 pts): Logistic Regression Model

- (2 pts) Split the data into features (X) and target (y) variables. The target variable is the 'sold\_within\_3\_months' column. The 'Price' column should be included as a feature.

```
In [87]: y_df = data_df['sold_within_3_months']
X_df = data_df.drop(columns=['sold_within_3_months'])
```

- (2 pts) Then split the data into train test sets using a 80-20 split. Use `random_state=42` for reproducibility.

```
In [88]: # X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)
X_train, X_test, y_train, y_test = train_test_split(X_df, y_df, test_size=0.2, random_state=42)
print(f"No. of training samples: {len(X_train)} No. of testing samples: {len(X_test)}")
```

No. of training samples: 1148 No. of testing samples: 288

- (2 pts) Standardize the features using `StandardScaler` from `sklearn.preprocessing` and then add a constant column using `sm.add_constant()`.

```
In [89]: #
scaler = StandardScaler()
X_train = scaler.fit_transform(X_train)
X_test = scaler.transform(X_test)
X_train = sm.add_constant(X_train)
X_test = sm.add_constant(X_test)
```

- (2 pts) Fit a logistic regression model on the training dataset.

```
In [90]: logistic_regression_model = sm.Logit(y_train, X_train).fit()
print(logistic_regression_model.summary())

Warning: Maximum number of iterations has been exceeded.
Current function value: 0.165914
Iterations: 35

=====
Logit Regression Results
=====
Dep. Variable:      sold_within_3_months      No. Observations:      1148
Model:              Logit      DF Residuals:          1134
Method:             MLE      DF Model:              13
Date:               Tue, 03 Dec 2024      Pseudo R-squ.:      0.7511
Time:               22:45:02      Log-Likelihood:      -196.47
converged:          False      LL-Null:              -765.10
Covariance Type:    nonrobust      LLR p-value:          1.448e-237
=====
              coef      std err      z      P>|z      [0.025      0.975]
-----
const          -0.3520    19.015    -0.019    0.985   -37.621    36.915
x1             -0.9840    0.082   -11.325    0.000   -11.713   -0.257
x2             -0.6112    0.319    -1.918    0.055   -1.236    0.013
x3             -0.2607    0.211    -1.235    0.217   -0.674    0.153
x4             -1.5443    0.654    -2.360    0.018   -2.827   -0.262
x5             1.1502    0.083     2.329    0.020    0.252    2.052
x6             0.1231    0.541     0.228    0.828   -0.937    1.184
x7             -2.0212    1.026    -1.970    0.049   -4.033   -0.010
x8             0.0730    0.589     0.125    0.900   -1.081    1.229
x9             0.1427    0.129     1.105    0.269   -0.110    0.396
x10            -0.1535    0.135    -1.134    0.257   -0.419    0.112
x11            -6.2822    5403.705   -0.001    0.999   -1.06e+04    1.06e+04
x12            -3.7054    3212.435   -0.001    0.999   -6295.062    6292.551
x13            -6.3885    5448.479   -0.001    0.999   -1.07e+04    1.07e+04
=====
```

Possibly complete quasi-separation: A fraction 0.20 of observations can be perfectly predicted. This might indicate that there is complete quasi-separation. In this case some parameters will not be identified.  
/Users/isaiah/miniconda3/envs/data\_hw2\_dryrun/lib/python3.10/site-packages/statsmodels/base/model.py:607: ConvergenceWarning: Maximum Likelihood optimization failed to converge. Check mle\_retvals  
warnings.warn("Maximum Likelihood optimization failed to "

- (2 pts) Evaluate the model on the test dataset using the accuracy score metric. Report the accuracy score.

```
In [ ]: y_pred_p_logistic = logistic_regression_model.predict(X_test)
y_test_pred_logistic = np.where(y_pred_p_logistic > 0.5, 1, 0)

print(accuracy_fn(y_test, y_test_pred_logistic))
# Accuracy between 0.91-0.95 gives full points

0.9385555555555556
```

- (2 pts) Calculate the precision, recall, and F1-score.

```
In [92]: precision = precision_fn(y_test, y_test_pred_logistic)
recall = recall_fn(y_test, y_test_pred_logistic)
f1 = f1_score_fn(y_test, y_test_pred_logistic)

print("Precision:", precision)
print("Recall:", recall)
print("F1-Score:", f1)

# Comment:
# Any difference within 1% is considered correct. Otherwise, 1 point should be deducted for each me
```

Precision: 0.9375  
Recall: 0.9482758620689655  
F1-Score: 0.9428571428571428

- (2 pts) Suppose that your company is running short on cash and needs to sell the cars quickly. How should you adjust the threshold for the logistic regression model to ensure that the company can sell the cars as quickly as possible?

- A. Increase the threshold
- B. Decrease the threshold

In a more general sense, how does the choice of threshold affect the precision and recall of the model?

Your Response:

- Decrease the threshold. Because this will make more predictions positive.
- Increase the threshold will increase the precision and decrease the recall.

- (6 pts) Use binary search to find the optimal threshold that maximizes the F1-score. Implement a binary search algorithm to find the threshold that maximizes the f1-score of the logistic regression model on the training set. The search interval should be between 0 and 1, and the stopping criterion is 10 iterations. What is the optimal threshold and what difference does the optimal threshold make in the F1-score?

```
In [ ]: def get_f1_score(y_test, y_pred_p, threshold):
    y_pred = np.where(y_pred_p > threshold, 1, 0)
    return f1_score_fn(y_test, y_pred)

def binary_search_threshold(y_test, y_pred_p, low, high, iterations):
    for i in range(iterations):
        mid = (low + high) / 2
        f1_score_low = get_f1_score(y_test, y_pred_p, mid)
        f1_score_high = get_f1_score(y_test, y_pred_p, high)
        if f1_score_low > f1_score_high:
            high = mid
        else:
            low = mid
    return mid

threshold = binary_search_threshold(y_test, y_pred_p_logistic, 0, 1, 10)
print("Threshold:", threshold)

y_test_pred_logistic = np.where(y_pred_p_logistic > threshold, 1, 0)

optimal_f1 = f1_score_fn(y_test, y_test_pred_logistic)
print("Optimal F1-Score:", optimal_f1)

# comment:
# This question is flawed. Any attempt to solve this question should be given full points.

Threshold: 0.4677734375
Optimal F1-Score: 0.80032467532467539417
```

### 3.3(18 pts) Decision Tree Model

Use a Decision Tree model from `sklearn` to predict whether a car will be sold within the first three months.

Follow these steps to complete the task:

- (2 pts) Train a Decision Tree Classifier to predict the target variable (`sold_within_3_months`). You can reuse the train and test sets from the previous section. Set `random_state=42` for reproducibility in `DecisionTreeClassifier`.

```
In [93]: from sklearn.tree import DecisionTreeClassifier

decision_tree_classifier = DecisionTreeClassifier(random_state=42)
decision_tree_classifier.fit(X_train, y_train)
```

```
Out[93]: DecisionTreeClassifier
DecisionTreeClassifier(random_state=42)
```

- (2 pts) Evaluate the model on the test set and report the depth of the tree.

```
In [ ]: # This recall looks not very promising, you wonder if you had a bug in your case.
# you decide to calculate a baseline which all predictions are randomly 0 or 1 with the same probab
y_test_pred_rf = decision_tree_classifier.predict(X_test)

acc = accuracy_fn(y_test, y_test_pred_rf)
depth = decision_tree_classifier.get_depth()
# comments:
# Any difference within 1% is considered correct. Otherwise, 1 point should be deducted for each me

print(f"Accuracy: {acc}, Depth: {depth}")

# 1% difference is acceptable, otherwise 1 point should be deducted for each metric, max 2 points c
# depth should be 16, if not, check what could be wrong, if issue comes from incorrect applicatio
```

Accuracy: 0.9201388888888889, Depth: 16

- (2 pts) Visualize the Decision Tree

```
In [95]: # Visualize the decision tree
from sklearn.tree import plot_tree
plt.figure(figsize=(20, 10))
plot_tree(decision_tree_classifier, filled=True)
plt.show()
```

- (2 pts) Retrain the Decision Tree Classifier with a maximum depth of 8 and evaluate it on the test set. Compare and explain the results.

```
In [96]: # train a decision tree with max_depth=8
decision_tree_classifier_depth_8 = DecisionTreeClassifier(max_depth=8, random_state=42)
decision_tree_classifier_depth_8.fit(X_train, y_train)

y_test_pred_rf_depth_8 = decision_tree_classifier_depth_8.predict(X_test)
print(accuracy_fn(y_test, y_test_pred_rf_depth_8))

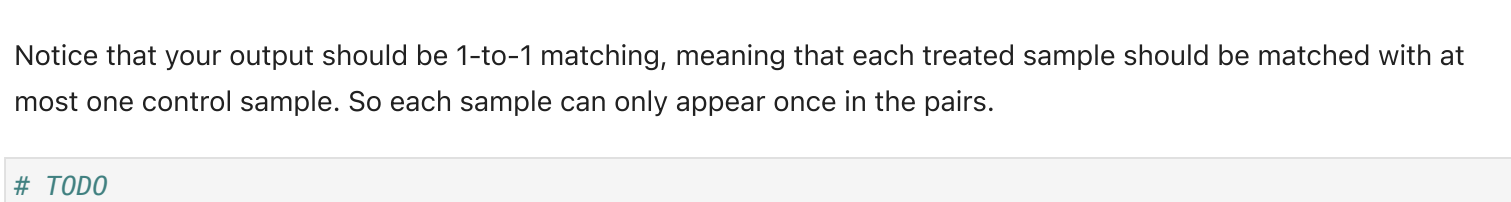
# comments:
# the tree with max_depth=8 should have a better performance than the one with the default depth=16
# you decide to calculate a baseline which all predictions are randomly 0 or 1 with the same probab
0.9236111111111112
```

- (6 pts) Train a Decision Tree Classifier for each depth from 1 to D where D is the maximum depth of the Decision Tree Classifier seen in the previous step. Evaluate each model on the test set and plot the accuracy of the models as a function of the depth and find the optimal depth.

```
In [99]: depths = []
accuracies = range(1, depth+1)
for d in depths:
    decision_tree_classifier = DecisionTreeClassifier(max_depth=d, random_state=42)
    decision_tree_classifier.fit(X_train, y_train)
    y_test_pred_rf = decision_tree_classifier.predict(X_test)
    acc = accuracy_fn(y_test, y_test_pred_rf)
    accuracies.append(acc)

plt.plot(depths, accuracies)
plt.xlabel('Depth')
plt.ylabel('Accuracy')
plt.title('Accuracy vs Depth')
plt.show()
```

# The optimal depth is 1, which is very counter-intuitive. You decide to investigate further by plotting accuracy vs depth.



- (4 pts) Train a decision tree of depth = 1, visualize the tree and explain what is the decision rule at the root node.

```
In [ ]: decision_tree_classifier_depth_1 = DecisionTreeClassifier(max_depth=1, random_state=42)
decision_tree_classifier_depth_1.fit(X_train, y_train)

plt.figure(figsize=(20, 10))
plot_tree(decision_tree_classifier_depth_1, filled=True)
plt.show()
```

# comment: Price is used as the first split of the tree, if it is lower than -0.036, the car is classified as sold within 3 months.

x[1] <= -0.036  
gini = 0.474  
samples = 1148  
value = [442, 706]

gini = 0.115  
samples = 717  
value = [44, 673]

gini = 0.141  
samples = 431  
value = [398, 33]

## Part 4 Propensity Score Matching (10 pts)

Your mentor is thrilled with the progress, and he has asked Ivan to put the model into production. Based on the model's prediction, the sales manager Ivan will decide whether to lower the car's price by 5%.

A new quarter has passed, and Ivan has collected updated sales data, which includes the following columns:

- Price: The initial price of the car.
- Pred\_Prob: The predicted probability of the car being sold within the first three months.
- Applied\_Discount: Whether the discount was applied (Yes=1, No=0).
- Discounted\_Price: The car's final price, calculated as `Price * 0.95` if the discount was applied; otherwise, it's equal to Price.
- Sold\_within\_3\_months: Whether the car was sold within the first three months (Yes=1, No=0).

Your task is to estimate the causal effect of the discount on sales within the first three months using propensity score matching.

```
In [100]: data_df = pd.read_csv('data/Task4.ToyotaCorolla_discount_sales.csv', index_col=0)
print(data_df.head())

   Price  Pred_Prob  Applied_Discount  Discounted_Price  Sold_within_3_months
0  12750         0.01             1         11475.0           1
1  21950         0.00             1         19755.0           1
2   9950         0.79             0          9950.0           1
3   9930         0.91             1         8937.0           1
4   9450         0.97             0          9450.0           0
```

- (1 pts) How many samples are in the treated group, and how many are in the control group?

```
In [105]: print(data_df['Applied_Discount'].value_counts())

# 118 are in the treated group and 82 are in the control group
Applied_Discount
0    118
1     82
Name: count, dtype: int64
```

### 4.2 (5 pts): Propensity Score Matching

- The propensity score is the predicted probability of the car being sold within the first three months from the logistic regression model, i.e. `Pred_Prob` column in the `Task4.ToyotaCorolla_discount_sales.csv`.

- Create pairs of matched samples as follows:
  - For each treated sample (discount applied), find a control sample (discount not applied) with a difference in propensity score of less than 0.05.
  - If there is more than one control sample for a treated sample, choose the control sample with the smallest difference in propensity score.
  - If there is no control sample satisfying the condition, discard the treated sample.

Notice that your output should be 1-to-1 matching, meaning that each treated sample should be matched with at most one control sample. So each sample can only appear once in the pairs.

```
In [118]: # 7000
treatment_df = data_df[data_df['Applied_Discount'] == 1]
control_df = data_df[data_df['Applied_Discount'] == 0]

print(f"there are {treatment_df.shape[0]} samples in the treated group.")
print(f"there are {control_df.shape[0]} samples in the control group.")
```

```
import networkx as nx
def get_similarity(propensity_score1, propensity_score2):
    """Calculate similarity for instances with given propensity scores"""
    return 1 - np.abs(propensity_score1 - propensity_score2)
```

```
G = nx.Graph()

for treatment_id, treatment_row in treatment_df.iterrows():
    for control_id, control_row in control_df.iterrows():
        ordered_matching.append((control_id, treatment_id))
        if similarity>0.95:
            G.add_weighted_edges_from([(control_id, treatment_id, similarity)])
```

```
matching = nx.max_weight_matching(G)
print(f"We have {len(matching)} successful matches!")
```

There are 82 samples in the treated group.  
There are 118 samples in the control group.  
We have 49 successful matches!

```
In [ ]: # 4.3 (4 pts): Average Treatment Effect (ATE)

Now let's estimate the effect of the discount on sales.
```

For each matched pair, there is one treated sample and one control sample. They may have different outcomes and we can calculate the average treatment effect (ATE) as

$$ATE = \frac{1}{N} \sum_{i=1}^N y_{treat}^{(i)} - y_{control}^{(i)}$$

where  $y_{treat}^{(i)}$  and  $y_{control}^{(i)}$  are the outcomes for the treated and control samples, respectively.

Notice that here the outcome is a simple binary variable, which is whether the car was sold within the first three months or not.

- (3 pts) Calculate the ATE based on the matched pairs and report the result

```
In [122]: ordered_matching = []

for control_id, treatment_id in matching:
    if control_id in control_df.index and treatment_id in treatment_df.index:
        ordered_matching.append((control_id, treatment_id))
    elif treatment_id in control_df.index and control_id in treatment_df.index:
        ordered_matching.append((treatment_id, control_id))
```

```
ATE = 0
for control_id, treatment_id in ordered_matching:
    treatment_value = treatment_df.loc[treatment_id, "Sold_within_3_months"]
    control_value = control_df.loc[control_id, "Sold_within_3_months"]
    ATE += treatment_value - control_value
```

ATE / # len(ordered\_matching)

```
print(f"Value of ATE: {ATE:.4f}")

# comment:
# In the matching pairs, sometimes the treatment is the first element, sometimes the control is the
# if the student didn't consider this, they will get a value lower than 0.429, like 0.388. 3 pts s
```

- (1 pts) What is your conclusion about the effect of the discount on sales within the first three months?

Your Response:

- The discount has a clear positive effect on sales within the first three months.