

P1 (50pt): Write a Python code in Colab using NumPy, Panda, Scikit-Learn to complete the following tasks:

```
import pandas as pd
import numpy as np
np.random.seed(100)
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression
from sklearn.metrics import mean_squared_error
from sklearn import preprocessing
from sklearn.preprocessing import StandardScaler
from sklearn.ensemble import RandomForestClassifier
from sklearn.model_selection import cross_val_score
from sklearn.model_selection import GridSearchCV
```

1. Import the Auto MPG dataset with `pandas.read_csv()` using the dataset URL, use the attribute names as explained in the dataset description as the column names (**5pt**), view the strings '?' as the missing value, and whitespace (i.e., '\s+') as the column delimiter. Print out the shape and first 5 rows of the obtained DataFrame. (**5pt**)

- Dataset source file: <http://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data>
- Dataset description: <http://archive.ics.uci.edu/ml/datasets/Auto+MPG>
- To have reproducible results using Scikit-Learn in Jupyter Notebook for the training/testing, set the seed at the beginning of your notebook https://www.mikulskibartos.name/how-to-set-the-global-random_state-in-scikit-learn/

```
# Put your answers here
url = "http://archive.ics.uci.edu/ml/machine-learning-databases/auto-mpg/auto-mpg.data"
df = pd.read_csv(url, index_col = False, names = ['mpg', 'cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model year',
        header = None, na_values = '?', sep = '\s+')

#df.columns = ['mpg', 'cylinders', 'displacement', 'horsepower', 'weight', 'acceleration', 'model year', 'origin', 'car name']

df.shape

(398, 9)

df.head(5)
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin	car name
0	18.0	8	307.0	130.0	3504.0	12.0	70	1	chevrolet chevelle malibu
1	15.0	8	350.0	165.0	3693.0	11.5	70	1	buick skylark 320
2	18.0	8	318.0	150.0	3436.0	11.0	70	1	plymouth satellite
3	16.0	8	304.0	150.0	3433.0	12.0	70	1	amc rebel sst
4	17.0	8	302.0	140.0	3449.0	10.5	70	1	ford torino

2. Delete the “car_name” column using .drop() method as it is irrelevant to the prediction. Print out a concise summary of the new DataFrame using .info() and check if NULL value exists in each column **(5pt)**

```
# Put your answers here
df = df.drop(columns = 'car name')
df.info()

<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0    mpg             398 non-null   float64
1    cylinders        398 non-null   int64
2    displacement     398 non-null   float64
3    horsepower       392 non-null   float64
4    weight           398 non-null   float64
5    acceleration     398 non-null   float64
6    model year       398 non-null   int64
7    origin           398 non-null   int64
dtypes: float64(5), int64(3)
memory usage: 25.0 KB
```

```
df.isnull()
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin
0	False	False	False	False	False	False	False	False
1	False	False	False	False	False	False	False	False
2	False	False	False	False	False	False	False	False
3	False	False	False	False	False	False	False	False
4	False	False	False	False	False	False	False	False
...
393	False	False	False	False	False	False	False	False
...

```
df.isnull().sum()
```

```
mpg          0
cylinders    0
displacement 0
horsepower   6
weight       0
acceleration 0
model year   0
origin       0
dtype: int64
```

3. Replace the NULL value with the mean value of the column using `.fillna()`. Print out the concise summary of the new DataFrame and recheck if NULL value exists in each column (5pt)

```
# Put your answers here
df['horsepower'].fillna(float(df['horsepower'].mean()), inplace = True)
df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 398 entries, 0 to 397
Data columns (total 8 columns):
#   Column          Non-Null Count  Dtype
---  -
0   mpg             398 non-null   float64
1   cylinders        398 non-null   int64
2   displacement     398 non-null   float64
3   horsepower       398 non-null   float64
4   weight           398 non-null   float64
5   acceleration     398 non-null   float64
6   model year       398 non-null   int64
7   origin           398 non-null   int64
```

```
dtypes: float64(5), int64(3)
memory usage: 25.0 KB
```

```
df.isnull().sum()
```

```
mpg          0
cylinders     0
displacement  0
horsepower    0
weight        0
acceleration  0
model year    0
origin        0
dtype: int64
```

4. For the 'origin' column with categorical attribute, replace it with the columns with numerical attributes using one-hot encoding (you can use either `get_dummies()` in Pandas or `OneHotEncoder` in Scikit-Learn). Print out the first 5 rows of the newly obtained DataFrame. **(10pt)**

a. <https://stackabuse.com/one-hot-encoding-in-python-with-pandas-and-scikit-learn/>

```
# Put your answers here
```

```
df = pd.get_dummies(data= df, columns=['origin'])
df.head(5)
```

	mpg	cylinders	displacement	horsepower	weight	acceleration	model year	origin_1	origin_2	origin_3
0	18.0	8	307.0	130.0	3504.0	12.0	70	1	0	0
1	15.0	8	350.0	165.0	3693.0	11.5	70	1	0	0
2	18.0	8	318.0	150.0	3436.0	11.0	70	1	0	0
3	16.0	8	304.0	150.0	3433.0	12.0	70	1	0	0
4	17.0	8	302.0	140.0	3449.0	10.5	70	1	0	0

5. Learn a linear regression model (`fit_intercept=True`) to predict the "mpg" column from the remaining columns in the obtained DataFrame of Step 4.

a. Separate the "mpg" column from other columns and view it as the label vector and others as the feature matrix **(5pt)**

```
# Put your answers here
```

```
X = df.drop('mpg', axis = 1)
```

```
y = df['mpg']
```

b. Split the data into a training set (80%) and testing set (20%) using `train_test_split` and print out their shapes **(5pt)**

```
# Put your answers here
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 50, train_size = 0.8)
```

```
X_train.shape
```

```
(318, 9)
```

```
X_test.shape
```

```
(80, 9)
```

```
y_train.shape
```

```
(318,)
```

```
y_test.shape
```

```
(80,)
```

c. Train the model using the training set and print out the coefficients of the model **(5 pt)**

```
# Put your answers here
```

```
model = LinearRegression()
```

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

```
print(model.coef_)
```

```
[-0.52953834  0.02398839 -0.0057078  -0.00698204  0.14021283  0.78677422  
 -1.92213942  0.7771369   1.14500252]
```

```
pd.DataFrame(model.coef_, X_train.columns, columns = ['Coeff'])
```

	Coeff
cylinders	-0.529538
displacement	0.023988
horsepower	-0.005708
weight	-0.006982
acceleration	0.140213
model year	0.786774

d. Use the learned model to predict on the test set and print out the mean squared error of the predictions **(5pt)**

```

origin 2      0.777137
# Put your answers here

pred = model.predict(X_test)
mean_squared_error(y_test, pred)

10.644550575766496

```

P2 (50pt): Write a Python code in Colab using NumPy, Panda, Scikit-Learn to complete the following tasks:

1. Import the red wine dataset with `pandas.read_csv()` using the dataset URL, use the semi-colon as the column delimiter, and print out both the first five rows and a concise summary of the obtained DataFrame. **(10 pt)**
 - a. Dataset source file: <http://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv>
 - b. Dataset description: <http://archive.ics.uci.edu/ml/datasets/wine+quality>
 - c. To have reproducible results using Scikit-Learn in Jupyter Notebook for the training/testing, set the seed at the beginning of your notebook https://www.mikulskibartosz.name/how-to-set-the-global-random_state-in-scikit-learn/

```

# Put your answers here
url2 = "http://archive.ics.uci.edu/ml/machine-learning-databases/wine-quality/winequality-red.csv"
df2 = pd.read_csv(url2, index_col = False, sep=';')
df2.head(5)

```

	fixed acidity	volatile acidity	citric acid	residual sugar	chlorides	free sulfur dioxide	total sulfur dioxide	density	pH	sulphates	alcohol	quality
0	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5
1	7.8	0.88	0.00	2.6	0.098	25.0	67.0	0.9968	3.20	0.68	9.8	5
2	7.8	0.76	0.04	2.3	0.092	15.0	54.0	0.9970	3.26	0.65	9.8	5
3	11.2	0.28	0.56	1.9	0.075	17.0	60.0	0.9980	3.16	0.58	9.8	6
4	7.4	0.70	0.00	1.9	0.076	11.0	34.0	0.9978	3.51	0.56	9.4	5

```
df2.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 1599 entries, 0 to 1598
Data columns (total 12 columns):
#   Column              Non-Null Count  Dtype
---  -
0   fixed acidity       1599 non-null   float64
1   volatile acidity    1599 non-null   float64
2   citric acid         1599 non-null   float64
3   residual sugar      1599 non-null   float64
4   chlorides           1599 non-null   float64
5   free sulfur dioxide 1599 non-null   float64
6   total sulfur dioxide 1599 non-null   float64
7   density             1599 non-null   float64
8   pH                  1599 non-null   float64
9   sulphates           1599 non-null   float64
10  alcohol             1599 non-null   float64
11  quality              1599 non-null   int64
dtypes: float64(11), int64(1)
memory usage: 150.0 KB
```

2. Suppose we want to predict the quality of wine from other attributes. Divide the data into a label vector and a feature matrix. Then split them into a training set (80%) and testing set (20%) using `train_test_split`. **(5pt)**

```
# Put your answers here
X = df2.drop('quality', axis = 1)
y = df2['quality']

X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=50, train_size=0.8)
```

3. Use the `StandardScaler()` in Scikit-Learn to preprocess the feature matrices of both training set and testing set. Note that the testing set can only be scaled by the mean and standard deviation values obtained from the training set. **(5 pt)**

a. <https://scikit-learn.org/stable/modules/preprocessing.html> (Refer to the code examples when using MinMaxScaler() in Section 6.3.1.1)

```
# Put your answers here
sc = StandardScaler()
X_train = sc.fit_transform(X_train)
X_test = sc.fit_transform(X_test)
```

4. Use the preprocessed training dataset to learn a classification model with RandomForestClassifier (with 300 trees) in Scikit-Learn. Use 5-fold cross-validation to train and cross-validate the model on the preprocessed training dataset. Print out the accuracies returned by the five folds as well as their average and standard deviation values. **(10 pt)**

```
# Put your answers here
classifier = RandomForestClassifier(n_estimators=300, random_state=0)
all_accuracies = cross_val_score(estimator=classifier, X=X_train, y=y_train, cv=5)

print(all_accuracies)

[0.66796875  0.6484375  0.68359375  0.66796875  0.68235294]

print(all_accuracies.mean())

0.6700643382352942

print(all_accuracies.std())

0.01273230424631959
```

5. Use GridSearchCV() to find and print out the best hyperparameter for the number of trees (try the following values: 100, 300, 500, 800, 1000 for 'n_estimator'), ignoring the search for other hyperparameters. Also use 5-fold cross-validation during GridSearchCV. **(15 pt)**

```
# Put your answers here
grid_param = {'n_estimators': [100, 300, 500, 800, 1000],
              'criterion': ['gini', 'entropy'],
              'bootstrap': [True, False]}

gd_sr = GridSearchCV(estimator=classifier,
                     param_grid=grid_param,
                     scoring='accuracy',
                     cv = 5,
                     n_jobs=-1)
```



```
gd_sr.fit(X_train, y_train)

GridSearchCV(cv=5,
             estimator=RandomForestClassifier(n_estimators=300, random_state=0),
             n_jobs=-1,
             param_grid={'bootstrap': [True, False],
                          'criterion': ['gini', 'entropy'],
                          'n_estimators': [100, 300, 500, 800, 1000]},
             scoring='accuracy')

best_parameters = gd_sr.best_params_
print(best_parameters)

{'bootstrap': True, 'criterion': 'entropy', 'n_estimators': 500}
```

6. Find and print out the testing accuracy of the model obtained using the best hyperparameter value on the preprocessed testing dataset (5 pt)

```
# Put your answers here
gd_sr.fit(X_train, y_train)

print(f"best mean cross-validation score: {gd_sr.best_score_}")
print(f"best parameters: {gd_sr.best_params_}")
print(f"test-set score: {gd_sr.score(X_test, y_test):.3f}")

best mean cross-validation score: 0.6755422794117647
best parameters: {'bootstrap': True, 'criterion': 'entropy', 'n_estimators': 500}
test-set score: 0.700
```

[Colab paid products](#) - [Cancel contracts here](#)

