

In-Lab

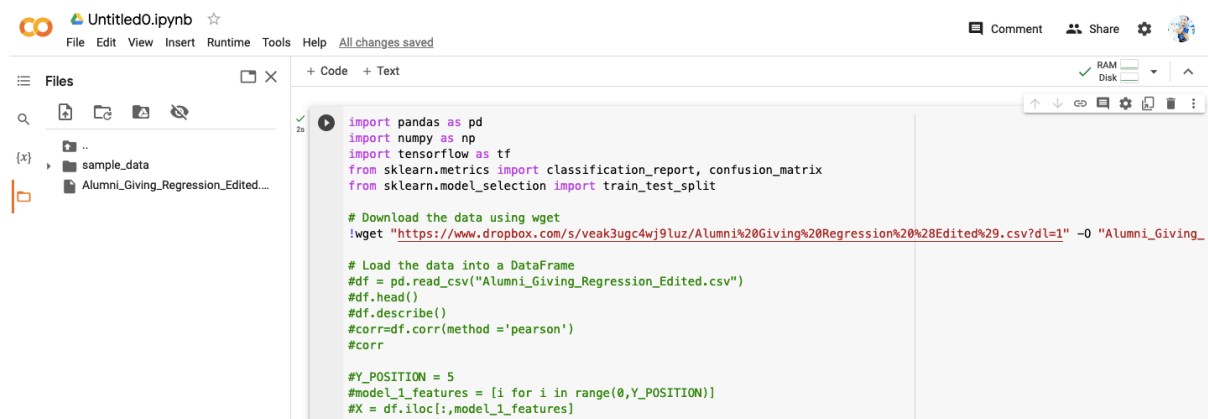
Task 1

```
import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
!wget
"https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"
```

Output:

We have successfully downloaded the dataset from link as you can see it in figure.



Task 2

```
import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget
"https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
```

```
df.head()
```

Output:

We are Loading and displaying the dataset using pandas. The `pd.read_csv("Alumni_Giving_Regression_Edited.csv")` function reads the dataset into a DataFrame, and `df.head()` displays the first few rows of the DataFrame (By Default show 5 rows).

The screenshot shows a Jupyter Notebook interface. On the left, the 'Files' pane shows a directory structure with 'sample_data' and 'Alumni_Giving_Regression_Edited...'. The main area displays the following code:

```
import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget "https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
df.head()
```

Below the code, the output of `df.head()` is displayed as a table:

	A	B	C	D	E	F
0	24	0.42	0.16	0.59	0.81	0.08
1	19	0.49	0.04	0.37	0.69	0.11
2	18	0.24	0.17	0.66	0.87	0.31
3	8	0.74	0.00	0.81	0.88	0.11
4	8	0.95	0.00	0.86	0.92	0.28

Task 3

```
import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget
"https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
df.describe()
```

Output:

We are summarizing the dataset statistics. The `df.describe()` function provides summary statistics (count, mean, std, min, max, etc.) for each numerical column in the DataFrame.

```

import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget "https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
df.describe()

```

	A	B	C	D	E	F
count	123.000000	123.000000	123.000000	123.000000	123.000000	123.000000
mean	17.772358	0.403659	0.136260	0.645203	0.841138	0.141789
std	4.517385	0.133887	0.060101	0.168794	0.083942	0.080674
min	6.000000	0.140000	0.000000	0.260000	0.580000	0.020000
25%	16.000000	0.320000	0.090000	0.500000	0.780000	0.080000
50%	18.000000	0.380000	0.130000	0.640000	0.840000	0.130000
75%	20.000000	0.480000	0.180000	0.780000	0.910000	0.170000
max	31.000000	0.950000	0.310000	0.960000	0.980000	0.410000

Task 4

```

import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget
"https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
corr=df.corr(method='pearson')
corr

```

Output:

We are calculating and displaying the Pearson correlation matrix for the dataset. The corr variable stores the correlation matrix calculated using `df.corr(method='pearson')`.

The screenshot shows a Jupyter Notebook titled 'Untitled0.ipynb'. The code in the cell is as follows:

```
import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget "https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
corr=df.corr(method='pearson')
corr
```

Below the code, a correlation matrix is displayed as a table:

	A	B	C	D	E	F
A	1.000000	-0.691900	0.414978	-0.604574	-0.521985	-0.549244
B	-0.691900	1.000000	-0.581516	0.487248	0.376735	0.540427
C	0.414978	-0.581516	1.000000	0.017023	0.055766	-0.175102
D	-0.604574	0.487248	0.017023	1.000000	0.934396	0.681660
E	-0.521985	0.376735	0.055766	0.934396	1.000000	0.647825
F	-0.549244	0.540427	-0.175102	0.681660	0.647825	1.000000

Task 5

```
import pandas as pd
import numpy as np
import tensorflow as tf
from sklearn.metrics import classification_report, confusion_matrix
from sklearn.model_selection import train_test_split

# Download the data using wget
#!wget
"https://www.dropbox.com/s/veak3ugc4wj9luz/Alumni%20Giving%20Regression%20%28Edited%29.csv?dl=1" -O "Alumni_Giving_Regression_Edited.csv"

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")
Y_POSITION = 5
model_1_features = [i for i in range(0,Y_POSITION)]
X = df.iloc[:,model_1_features]
Y = df.iloc[:,Y_POSITION]
X_train, X_test, y_train, y_test = train_test_split(X, Y,
test_size=0.20, random_state=2020)
```

Output:

We are splitting the dataset into training and testing sets for machine learning. The code uses `train_test_split` from scikit-learn to split the features (X) and the target variable (Y) into training and testing sets.

Task 6

```
import pandas as pd
import numpy as np
from sklearn.metrics import mean_squared_error
from sklearn.model_selection import train_test_split
from sklearn.linear_model import LinearRegression

# Load the data into a DataFrame
df = pd.read_csv("Alumni_Giving_Regression_Edited.csv")

Y_POSITION = 5
model_1_features = [i for i in range(0, Y_POSITION)]
X = df.iloc[:, model_1_features]
Y = df.iloc[:, Y_POSITION]

X_train, X_test, y_train, y_test = train_test_split(X, Y,
test_size=0.20, random_state=2020)

modell = LinearRegression()
modell.fit(X_train, y_train)
y_pred_train1 = modell.predict(X_train)

print("Regression")
print("=====")
RMSE_train1 = mean_squared_error(y_train, y_pred_train1)
print("Regression Train set: RMSE ", RMSE_train1)

print("=====")
y_pred1 = modell.predict(X_test)
RMSE_test1 = mean_squared_error(y_test, y_pred1)
print("Regression Test set: RMSE ".format(RMSE_test1))
print("=====")
coef_dict = {}
for coef, feat in zip(modell.coef_, model_1_features):
    coef_dict[df.columns[feat]] = coef
print(coef_dict)
```

Output:

We are building a linear regression model, fitting it to the training data, and evaluating the model's performance. The code uses scikit-learn to create a linear regression model, train it,

make predictions on the training and testing data, and calculate the Root Mean Squared Error (RMSE) for model evaluation. It also displays the coefficients of the model for each feature.

```

ision_Edited...
Y_POSITION = 5
model_1_features = [i for i in range(0, Y_POSITION)]
X = df.iloc[:, model_1_features]
Y = df.iloc[:, Y_POSITION]

X_train, X_test, y_train, y_test = train_test_split(X, Y, test_size=0.20, random_state=2020)

model1 = LinearRegression()
model1.fit(X_train, y_train)
y_pred_train1 = model1.predict(X_train)

print("Regression")
print("=====")
RMSE_train1 = mean_squared_error(y_train, y_pred_train1)
print("Regression Train set: RMSE ", RMSE_train1)

print("=====")
y_pred1 = model1.predict(X_test)
RMSE_test1 = mean_squared_error(y_test, y_pred1)
print("Regression Test set: RMSE ", format(RMSE_test1))
print("=====")
coef_dict = {}
for coef, feat in zip(model1.coef_, model_1_features):
    coef_dict[df.columns[feat]] = coef
print(coef_dict)

```

```

Regression
=====
Regression Train set: RMSE  0.002761693322289229
=====
Regression Test set: RMSE
=====
{'A': -0.0009337757382416938, 'B': 0.16012156890162943}
{'A': -0.0009337757382416938, 'B': 0.16012156890162943, 'C': -0.044160015425349614}
{'A': -0.0009337757382416938, 'B': 0.16012156890162943, 'C': -0.044160015425349614, 'D': 0.15217907817100407}
{'A': -0.0009337757382416938, 'B': 0.16012156890162943, 'C': -0.044160015425349614, 'D': 0.15217907817100407, 'E': 0.17539950794101047}

```

Note:

The tasks are performed using common libraries like pandas, numpy, scikit-learn (sklearn), and TensorFlow. However, we have flexibility in choosing alternative libraries or tools for similar tasks like Dask, Vaex, cuDF (for large datasets or distributed computing). We can use Matplotlib or Seaborn for more advanced data visualization. We can use NumPy for basic statistics, and Pandas Profiling for more detailed data profiling. We can use SciPy for correlation functions. We can use TensorFlow's built-in data preprocessing functions. We can use other machine learning libraries like XGBoost, LightGBM, or scikit-learn for regression models.