

HW 6

This assignment covers all fundamental concepts required for completing a project

DO NOT ERASE MARKDOWN CELLS AND INSTRUCTIONS IN YOUR HW submission

- **Q** - QUESTION
- **A** - Where to input your answer

Instructions

Keep the following in mind for all notebooks you develop:

- Structure your notebook.
- Use headings with meaningful levels in Markdown cells, and explain the questions each piece of code is to answer or the reason it is there.
- Make sure your notebook can always be rerun from top to bottom.
- Please start working on this assignment as soon as possible. If you are a beginner in Python this might take a long time. One of the objectives of this assignment is to help you learn python and scikit-learn package.
- See [README.md](#) for homework submission instructions

Related Tutorials

Refreshers

- [Intro to Machine Learning w scikit-learn](#)
- [A tutorial on statistical-learning for scientific data processing](#)

Classification Approaches

- [Logistic Regression with Sklearn](#)
- [KNN with sklearn](#)
- [Support Vector machine example](#)
- [SVC](#)
- [Bagging Classifier](#)
- [Gradient Boosting Classifier](#)

Modeling

- [Cross-validation](#)
- [Plot Confusion Matrix with Sklearn](#)
- [Confusion Matrix Display](#)

Import all required library

```
[1]: import pandas as pd
from sklearn.model_selection import train_test_split
import matplotlib.pyplot as plt
import numpy as np
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler, MaxAbsScaler
import json
import lightgbm as lgbm
from sklearn.decomposition import PCA
```

```

from sklearn.manifold import TSNE
import seaborn as sns
from imblearn.over_sampling import RandomOverSampler
from sklearn.ensemble import RandomForestClassifier

```

Data Processing

Q1 Get training data from the dataframe

1. Load `HW6_data.csv` from `data` folder into data frame
2. Print the head of the dataframe
3. Print the shape of the dataframe
4. Print the description of the dataframe
5. Check if the dataset has NULL values. (Show number of NULL values per column)
6. Assign `Cover_Type` values to Y
7. Assign rest of the column values to X

A1 Fill the cell blocks below, Create new cell as per your necessary

[50]: *#You can create or remove cells as per your need*

```

[39]: df = pd.read_csv('../data/HW6_data.csv')
print(df.head())
print(df.shape)
print(df.describe())
df.isnull().sum()

X= df.drop('Cover_Type', axis = 1)
Y= df['Cover_Type']

   Elevation Aspect Slope Horizontal_Distance_To_Hydrology \
0      3080.0    137    18.0                           166
1      2758.0     19     8.0                           551
2      2779.0    86     9.0                           43
3      2811.0   296     0.0                          287
4      2956.0   314    26.0                           71

   Vertical_Distance_To_Hydrology Horizontal_Distance_To_Roadways \
0                               1                           1009
1                             49                           1766
2                            -10                          3889
3                                4                           788
4                             22                          2910

   Hillshade_9am Hillshade_Noon Hillshade_3pm \
0        250.0       198       166
1        225.0       231       124
2        155.0       204       123
3        191.0       226       113
4        230.0       200        99

   Horizontal_Distance_To_Fire_Points ... Soil_Type32 Soil_Type33 \
0            3635.0 ...          0          0
1            1648.0 ...          0          0
2            364.0 ...          0          0
3            144.0 ...          0          0
4            743.0 ...          0          0

   Soil_Type34 Soil_Type35 Soil_Type36 Soil_Type37 Soil_Type38 \
0            0       0.0          0          0          0
1            0       0.0          0          0          0
2            0       0.0          0          0          1
3            0       0.0          0          0          0
4            0       0.0          0          0          1

```

```

4          0        0.0        0          0          1

Soil_Type39  Soil_Type40  Cover_Type
0            0        0.0        1
1            0        0.0        2
2            0        0.0        2
3            0        0.0        2
4            0        NaN        2

[5 rows x 55 columns]
(80000, 55)

      Elevation      Aspect      Slope \
count  79433.000000  80000.000000  79346.000000
mean   2981.436531  151.634175  15.092494
std    287.979705  109.945631  8.528153
min   1813.000000 -29.000000 -3.000000
25%  2762.000000  60.000000  9.000000
50%  2967.000000 122.000000 14.000000
75%  3217.000000 246.000000 20.000000
max   4271.000000 400.000000 61.000000

Horizontal_Distance_To_Hydrology  Vertical_Distance_To_Hydrology \
count                      80000.000000                  80000.000000
mean                     271.564212                  51.510737
std                      227.532197                  68.091489
min                     -43.000000                 -276.000000
25%                    111.000000                  4.000000
50%                    212.000000                  31.000000
75%                    361.000000                  78.000000
max                     1544.000000                 562.000000

Horizontal_Distance_To_Roadways  Hillshade_9am  Hillshade_Noon \
count                      80000.000000                  79200.000000  80000.000000
mean                     1770.080712                 211.786818  221.069125
std                      1318.661060                  30.822278  22.191030
min                     -238.000000                 10.000000  69.000000
25%                    821.000000                  198.000000  210.000000
50%                    1440.000000                 218.000000  224.000000
75%                    2366.000000                 233.000000  237.000000
max                     7604.000000                 293.000000  264.000000

Hillshade_3pm  Horizontal_Distance_To_Fire_Points ...  Soil_Type32 \
count          80000.000000                  78870.000000 ...  80000.000000
mean         140.711750                  1578.058615 ...  0.038150
std          43.859689                  1125.780446 ...  0.191559
min         -48.000000                 -218.000000 ...  0.000000
25%        115.000000                  782.000000 ...  0.000000
50%        142.000000                 1362.000000 ...  0.000000
75%        169.000000                 2082.000000 ...  0.000000
max         268.000000                 8011.000000 ...  1.000000

Soil_Type33  Soil_Type34  Soil_Type35  Soil_Type36  Soil_Type37 \
count  80000.000000  80000.000000  79720.000000  80000.000000  80000.000000
mean   0.037687  0.011838  0.015429  0.010812  0.012538
std    0.190441  0.108155  0.123252  0.103420  0.111268
min   0.000000  0.000000  0.000000  0.000000  0.000000
25%  0.000000  0.000000  0.000000  0.000000  0.000000
50%  0.000000  0.000000  0.000000  0.000000  0.000000
75%  0.000000  0.000000  0.000000  0.000000  0.000000
max   1.000000  1.000000  1.000000  1.000000  1.000000

Soil_Type38  Soil_Type39  Soil_Type40  Cover_Type
count  80000.000000  80000.000000  75000.000000  80000.000000
mean   0.040325  0.039163  0.030707  1.770725
std    0.196722  0.193983  0.172523  0.892577
min   0.000000  0.000000  0.000000  1.000000
25%  0.000000  0.000000  0.000000  1.000000
50%  0.000000  0.000000  0.000000  2.000000
75%  0.000000  0.000000  0.000000  2.000000
max   1.000000  1.000000  1.000000  7.000000

```

```
[8 rows x 55 columns]
```

Q2: Open-Ended Questions: Observe the range of all feature values and statistical information from the dataframe description above.

1. If the dataset has NULL values, Give proper justification about the methods you will use to replace NULL values for specific columns.
2. Do you think in our dataset normalization is required? -- Give proper justification based on your opinion.
3. What type of normalization/Scaling technique you would recommend for our dataset?

A2

Answer 1: I would use the mean or median to replace NULL values, depending on the distribution of each variable.

Answer 2: I think normalization is required due to the different scales of features which can effect the performance of distance based algorithms. It ensures that each feature is given equal weight in a predictive model.

Answer 3: I recommend using StandardScaler for features following a Gaussian distribution and MinMaxScaler for bounded features without a Gaussian distribution.

Q3:

1. Replace the null values with the best possible methods from your above observation
2. Use the above mentioned normalization technique on our HW_6 dataset.
3. Transform the X dataframe using chosen normalization technique.

Note: Make sure the scaled X has all column name same as X dataframe

A3 Fill the cell blocks below, Create new cell as per your necessary

```
[40]: # Replace NULLs  
# You can create or remove cells as per your need  
from sklearn.impute import SimpleImputer  
imputer = SimpleImputer(strategy='mean')  
df_imputed = pd.DataFrame(imputer.fit_transform(df), columns=df.columns)  
  
X= df_imputed.drop('Cover_Type', axis=1)
```

```
[41]: # Normalize data  
# You can create or remove cells as per your need  
scaler = StandardScaler()  
Scaled_X = pd.DataFrame(scaler.fit_transform(X), columns=X.columns)
```

```
[ ]:
```

Q4:

1. Check again and show if there is any null values left in our Scaled_X .
2. Print all unique values/ different class id from the Y data .

A4 Fill the cell blocks below, Create new cell as per your necessary

```
[ ]: # You can create or remove cells as per your need
```

```
[42]: print(type(Scaled_X))  
null_values = Scaled_X.isnull().sum()  
print(null_values[null_values > 0])  
unique_classes = Y.unique()  
print(unique_classes)  
  
<class 'pandas.core.frame.DataFrame'>  
Series([], dtype: int64)  
[1 2 3 7 6 4]
```

Part 1: I have a subset of whole data(N=20000) for Data Visualization

Part 1: Use a subset of whole data(N=20000) for Data Visualization

Data Subset Creation

1. First we are Selecting `N=20000` random rows from our original dataset which is `df` and create a new subset of data.
2. Using the below `rndperm` and selecting first N index from the `Scaled_X` and `Y`

```
[43]: np.random.seed(42)
rndperm = np.random.permutation(df.shape[0])
N = 20000
data_subset_x = Scaled_X.loc[rndperm[:N], :].copy()
data_subset_y = Y.loc[rndperm[:N]].copy()
```

Q5:

1. Use PCA and reduce the dimension of the `data_subset_x` into `3`.
2. Store the PCA result into `pca_result` variable
3. Add the results from the PCA into the `data_subset_x` as new columns. (Choose any meaningful names for the columns)

A5 Fill the below cells. Use extra cells as per your necessary

```
[44]: #You can create or remove cells as per your need
from sklearn.decomposition import PCA
pca = PCA(n_components=3)
pca_result = pca.fit_transform(data_subset_x)
```

```
[45]: data_subset_x['pca_one'] = pca_result[:,0]
data_subset_x['pca_two'] = pca_result[:,1]
data_subset_x['pca_three'] = pca_result[:,2]
```

Q6:

1. Use TSNE and reduce the dimension of the `data_subset_x` into `2`.
2. Store the TSNE result into `tsne_results` variable
3. Add the results from the T-SNE into the `data_subset_x` as new columns. (Choose any meaningful names for the columns)

Note:

1. You can use `from sklearn.manifold import TSNE` for TSNE initialization.
2. Give value of `n_components` as per the question.
3. Also use other parameters while TSNE initialization as, `verbose=1, perplexity=40, n_iter=300`

A6 Fill the below cells. Use extra cells as per your necessary

```
[46]: #You can create or remove cells as per your need
from sklearn.manifold import TSNE
tsne = TSNE(n_components=2, verbose=1, perplexity=40, n_iter=300)
tsne_results = tsne.fit_transform(data_subset_x)
```

```
[t-SNE] Computing 121 nearest neighbors...
[t-SNE] Indexed 20000 samples in 0.007s...
[t-SNE] Computed neighbors for 20000 samples in 0.359s...
[t-SNE] Computed conditional probabilities for sample 1000 / 20000
[t-SNE] Computed conditional probabilities for sample 2000 / 20000
[t-SNE] Computed conditional probabilities for sample 3000 / 20000
[t-SNE] Computed conditional probabilities for sample 4000 / 20000
[t-SNE] Computed conditional probabilities for sample 5000 / 20000
[t-SNE] Computed conditional probabilities for sample 6000 / 20000
[t-SNE] Computed conditional probabilities for sample 7000 / 20000
[t-SNE] Computed conditional probabilities for sample 8000 / 20000
[t-SNE] Computed conditional probabilities for sample 9000 / 20000
[t-SNE] Computed conditional probabilities for sample 10000 / 20000
[t-SNE] Computed conditional probabilities for sample 11000 / 20000
```

```
[t-SNE] Computed conditional probabilities for sample 12000 / 20000
[t-SNE] Computed conditional probabilities for sample 13000 / 20000
[t-SNE] Computed conditional probabilities for sample 14000 / 20000
[t-SNE] Computed conditional probabilities for sample 15000 / 20000
```

```
[47]: data_subset_x['tsne-2d-one'] = tsne_results[:,0]
data_subset_x['tsne-2d-two'] = tsne_results[:,1]
```

Q7:

1. Create a new dataframe with name `df_plot`
2. This dataframe will merge everything from `data_subset_x` and `data_subset_y`
3. We need to give a name for the `data_subset_y` column. Use `Cover_Type` as the name of the column

A7 Fill the below cells. Use extra cells as per your necessary

```
[48]: df_plot = data_subset_x.copy()
df_plot['Cover_Type'] = data_subset_y.values
```

Q8: Now we will plot all points from our dataframe `df_plot` Using the result from **PCA**

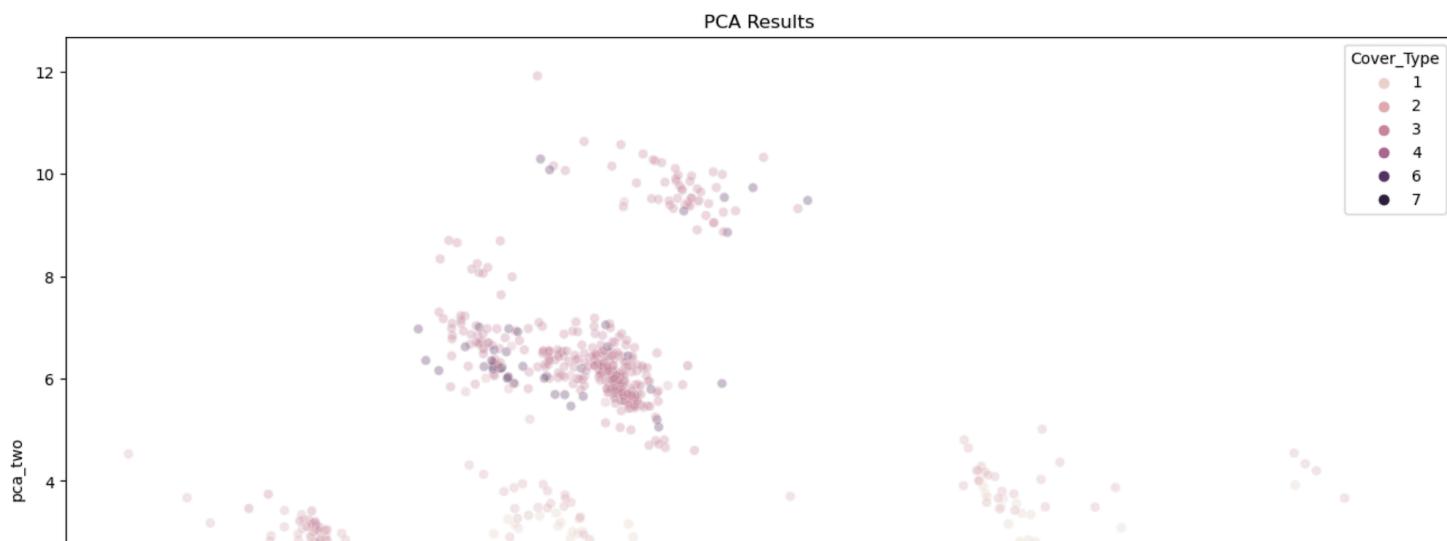
1. Use `pca-one` and `pca-two` column as X and Y axis respectively.
2. Use seaborn scatterplot for plotting the points.

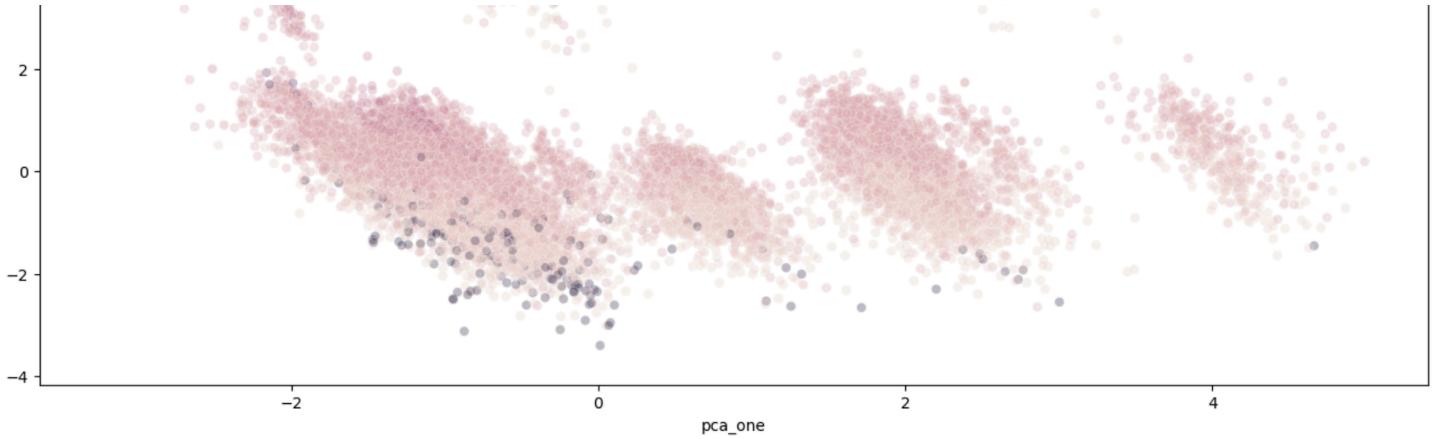
Note: Use the notebook from class for reference. The link is provided below.

Link: https://git.txstate.edu/ML/2022Fall/blob/main/project/Data_Viz_with_PCA_TSNE.ipynb

A8 Fill the below cells. Use extra cells as per your necessary

```
[49]: import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(16,10))
sns.scatterplot(
    x='pca_one', y='pca_two',
    hue='Cover_Type',
    data=df_plot,
    legend='full',
    alpha=0.3
)
plt.title('PCA Results')
plt.show()
```





Q9: Now we will plot all points from our dataframe `df_plot`. Using result from T-SNE.

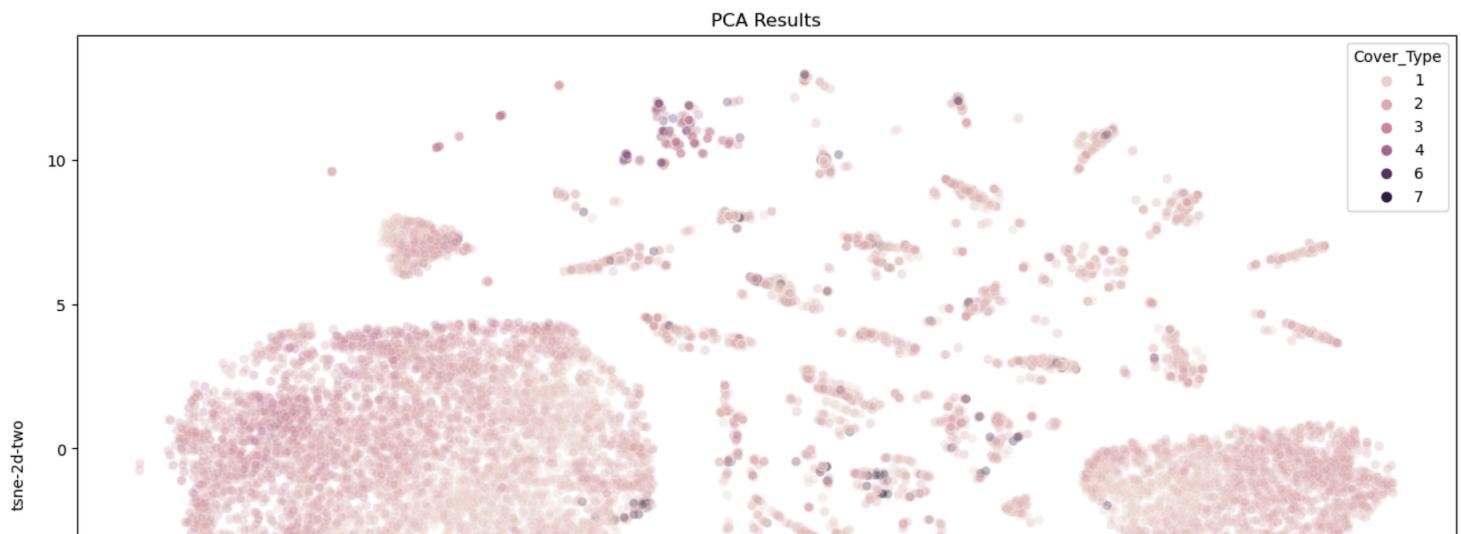
1. Use `tsne-2d-one` and `tsne-2d-two` column as X and Y axis respectively.
2. Use seaborn scatterplot for plotting the points.

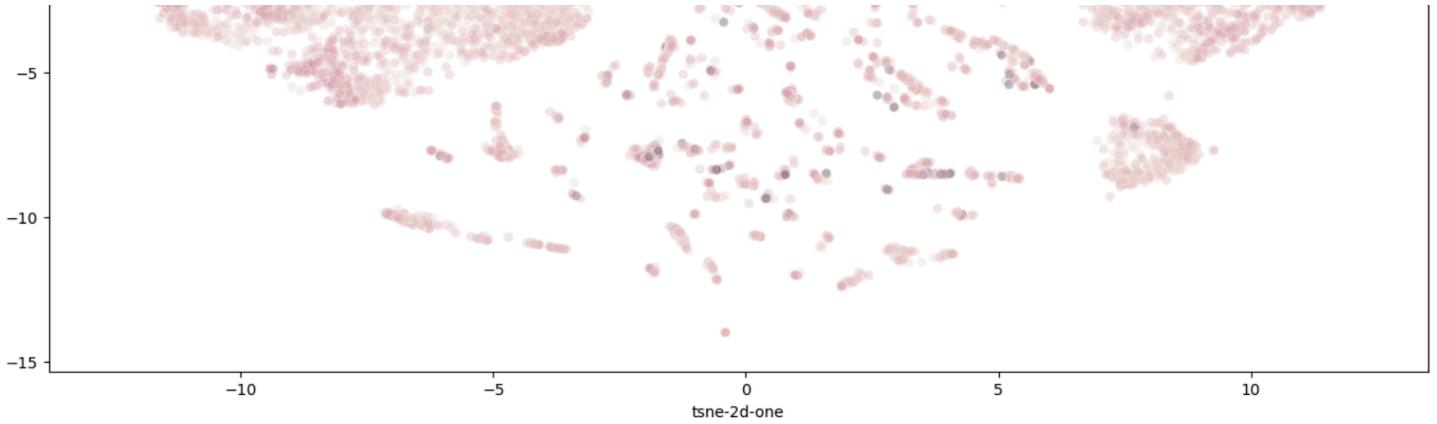
Note: Use the notebook from class for reference. The link is provided below.

Link: https://git.txstate.edu/ML/2022Fall/blob/main/project/Data_Viz_with_PCA_TSNE.ipynb

A9 Fill the below cells. Use extra cells as per your necessary

```
[50]: import matplotlib.pyplot as plt
import seaborn as sns
plt.figure(figsize=(16,10))
sns.scatterplot(
    x='tsne-2d-one', y='tsne-2d-two',
    hue='Cover_Type',
    data=df_plot,
    legend='full',
    alpha=0.3
)
plt.title('PCA Results')
plt.show()
```





Part 2: Data Analysis and Classification Using Entire Dataset

Q10: Observe the data plotting and find the relation between datapoints and their characteristics.

1. Reduce the dimension of our `Scaled_X` dataframe to 3 using PCA algorithm.
2. Store the result into a variable named `pca_result`
3. Create Train data and Test data using `pca_result` and Y.

Note:

1. Consider `pca_result` as X values, and Y as y values.
2. You can use sklearn `train_test_split`
3. Keep Train and Test ratio as : 75%:25%

A10 Fill the below cells. Use extra cells as per your necessary

```
[51]: #You can create or remove cells as per your need
pca = PCA(n_components=3)
pca_result = pca.fit_transform(Scaled_X)
```

```
[52]: x_train,x_test,y_train,y_test = train_test_split(
    pca_result, Y, test_size=0.25, random_state=42
)
```

Now, Select Three best model for our dataset. You have to decide three models which might work well with our dataset.

Q11

Model Number 1

1. Reason behind choosing the model.
2. Create the model using sklearn or any proper library
3. Fit the model with the train data
4. Get the score from the model using test data

A11 Fill the below cells. Use extra cells as per your necessary

Answer for Q.No:1 goes here The reason behind choosing LightGBM is due to its high efficiency and speed, which comes from its histogram-based learning.

```
[53]: from lightgbm import LGBMClassifier
from sklearn.metrics import accuracy_score
model = LGBMClassifier(random_state=42)
```



```
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
Accuracy of the LightGBM model: 0.7435
```

Q12

Model Number 2

1. Reason behind choosing the model.
2. Create the model using sklearn or any proper library
3. Fit the model with the train data
4. Get the score from the model using test data

A12 Fill the below cells. Use extra cells as per your necessaryReplace ??? with code in the code cell below

Answer for Q.No:1 goes here Random forests can capture really complex, non-linear relationships between features without requiring transformation.

```
[54]: rf_model = RandomForestClassifier(random_state=42)
rf_model.fit(x_train, y_train)
y_pred = rf_model.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy of the Random Forest model:", accuracy)
```

Accuracy of the Random Forest model: 0.7749

Q13

Model Number 3

1. Reason behind choosing the model.
2. Create the model using sklearn or any proper library
3. Fit the model with the train data
4. Get the score from the model using test data

A13 Fill the below cells. Use extra cells as per your necessary

Answer for Q.No:1 goes here KNN can capture non-linear decision boundaries without any extra effort, and no transformations are necessary for the input data.

```
[55]: from sklearn.neighbors import KNeighborsClassifier

knn_model = KNeighborsClassifier()

knn_model.fit(x_train, y_train)

y_pred = knn_model.predict(x_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy of the KNN model:", accuracy)
```

Accuracy of the KNN model: 0.76425

Q14

1. Plot a histogram using Y dataframe and display the per-class data distribution(number of rows per class).
2. Also print the number of rows per class as numeric value.

A14 Fill the below cells. Use extra cells as per your necessary

```
[ ]: #You can create or remove cells as per your need
```

```
[56]: class_counts = Y.value_counts()
print("Number of rows per class:")
print(class_counts)
plt.figure(figsize=(10, 6))
```

```

class_counts.plot(kind='bar')
plt.title('Per-Class Data Distribution')
plt.xlabel('Class')
plt.ylabel('Number of Rows')
plt.grid(True)
plt.show()

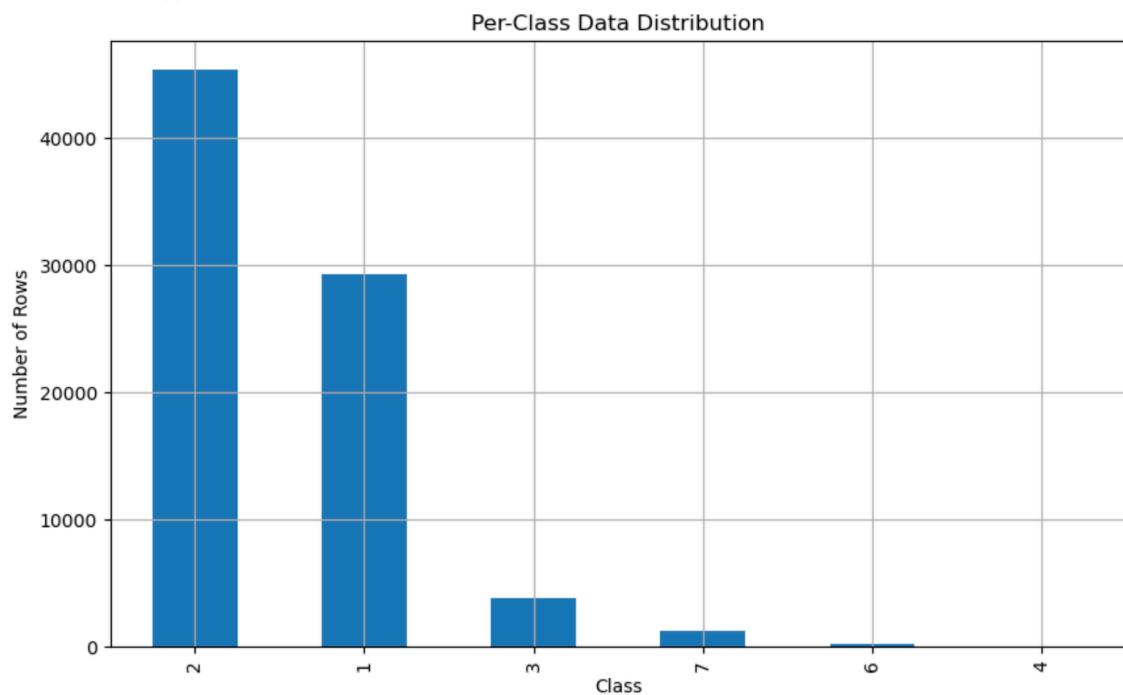
```

Number of rows per class:

Cover_Type

2	45393
1	29311
3	3816
7	1245
6	229
4	6

Name: count, dtype: int64



Q15

1. From the histogram we can see that the dataset is highly imbalanced.
2. Use a proper dataset balancing technique to make the dataset balanced.
3. Plot a histogram using new y values and display the per-class data distribution(number of rows per class).

Note: You can use the `imblearn.over_sampling` library for this task. But use appropriate strategy for the method.

Follow the documentation for details: https://imbalanced-learn.org/stable/references/generated/imblearn.over_sampling.SMOTE.html

A15 Fill the below cells. Use extra cells as per your necessary

[57]: #You can create or remove cells as per your need

```

ros = RandomOverSampler(random_state=42)
X_res, y_res = ros.fit_resample(x_train, y_train)

resampled_class_counts = pd.Series(y_res).value_counts()
print("Number of rows per class after resampling:")

```

```

print(resampled_class_counts)

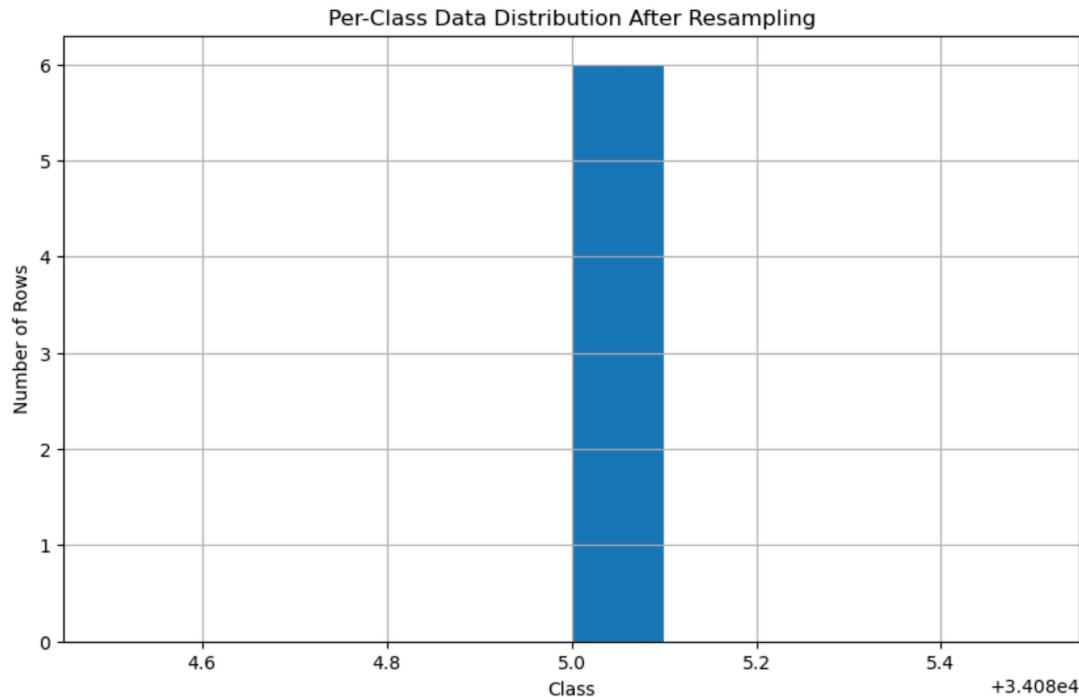
plt.figure(figsize=(10, 6))
resampled_class_counts.plot(kind='hist')
plt.title('Per-Class Data Distribution After Resampling')
plt.xlabel('Class')
plt.ylabel('Number of Rows')
plt.grid(True)
plt.show()

```

Number of rows per class after resampling:

Cover_Type	Count
1	34085
2	34085
3	34085
6	34085
7	34085
4	34085

Name: count, dtype: int64



[]:

Q16

1. Create new Train and Test data from the balanced X and Y value.
2. Keep Train and Test ratio as : 75%:25%

A16 Fill the below cells. Use extra cells as per your necessary

```

[58]: x_train,x_test,y_train,y_test = train_test_split(
    X_res, y_res, test_size=0.25, random_state=42
)

```

Q17

Now Use the previously initialized three models and calculate the score from our new balanced dataset

Now, use the previously mentioned three models and calculate the score from our new balanced dataset.

Model Number 1

1. Fit the model with the new train data(Use the previous Model 1)
2. Get the score from the model using new test data

A17 Fill the below cells. Use extra cells as per your necessary

```
[59]: #You can create or remove cells as per your need
model.fit(x_train, y_train)

y_pred = model.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy of the LightGBM model on balanced data:", accuracy)

[LightGBM] [Info] Auto-choosing row-wise multi-threading, the overhead of testing was 0.000658 seconds.
You can set `force_row_wise=true` to remove the overhead.
And if memory is not enough, you can set `force_col_wise=true`.
[LightGBM] [Info] Total Bins 765
[LightGBM] [Info] Number of data points in the train set: 153382, number of used features: 3
[LightGBM] [Info] Start training from score -1.787609
[LightGBM] [Info] Start training from score -1.788388
[LightGBM] [Info] Start training from score -1.792725
[LightGBM] [Info] Start training from score -1.787764
[LightGBM] [Info] Start training from score -1.798025
[LightGBM] [Info] Start training from score -1.796098
[LightGBM] [Warning] No further splits with positive gain, best gain: -inf
Accuracy of the LightGBM model on balanced data: 0.8818651228289782
```

Model Number 2

1. Fit the model with the new train data(Use the previous Model 2)
2. Get the score from the model using new test data

Fill the below cells. Use extra cells as per your necessary

```
[60]: #You can create or remove cells as per your need
rf_model.fit(x_train, y_train)

y_pred = rf_model.predict(x_test)

accuracy = accuracy_score(y_test, y_pred)
print("Accuracy of the Random Forest model on balanced data:", accuracy)

Accuracy of the Random Forest model on balanced data: 0.9498904709748083
```

Model Number 3

1. Fit the model with the new train data(Use the previous Model 3)
2. Get the score from the model using new test data

Fill the below cells. Use extra cells as per your necessary

```
[61]: #You can create or remove cells as per your need
knn_model.fit(x_train, y_train)

y_pred = knn_model.predict(x_test)
accuracy = accuracy_score(y_test, y_pred)
print("Accuracy of the KNN model:", accuracy)

Accuracy of the KNN model: 0.9167188233453294
```

After making the dataset balanced we can see a significant improve in the performance for all three models.