


Import the Libraries

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
```

Import dataset(csv file)

```
1 titanic_data = pd.read_csv("/content/train.csv")
2 titanic_data
```

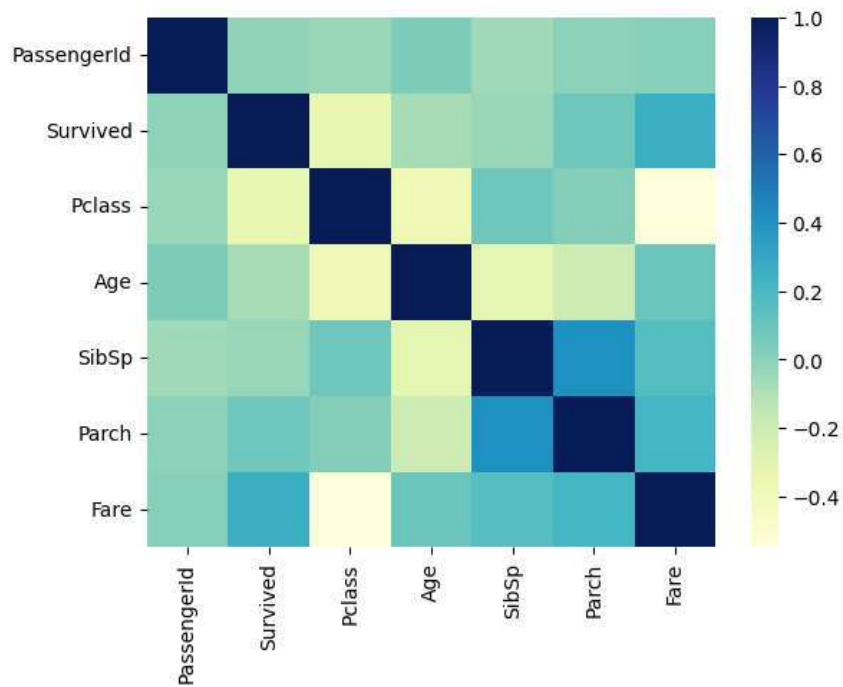


	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	1	0	3	Braund, Mr. Owen Harris	male	22.0	1	0	A/5 21171	7.2500	NaN	S
1	2	1	1	Cumings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikkinen, Miss. Laina	female	26.0	0	0	STON/O2. 3101282	7.9250	NaN	S
3	4	1	1	Futrelle, Mrs. Jacques Heath (Lily May Peel)	female	35.0	1	0	113803	53.1000	C123	S
4	5	0	3	Allen, Mr. William Henry	male	35.0	0	0	373450	8.0500	NaN	S
...
886	887	0	2	Montvila, Rev. Juozas	male	27.0	0	0	211536	13.0000	NaN	S
887	888	1	1	Graham, Miss. Margaret Edith	female	19.0	0	0	112053	30.0000	B42	S
888	889	0	3	Johnston, Miss. Catherine Helen "Carrie"	female	NaN	1	2	W./C. 6607	23.4500	NaN	S
889	890	1	1	Behr, Mr. Karl Howell	male	26.0	0	0	111369	30.0000	C148	C

Corrlation Heatmap

```
1 import seaborn as sns
2
3 sns.heatmap(titanic_data.corr(), cmap="YlGnBu")
4 plt.show()
```

<ipython-input-39-ea442d09ae1a>:3: FutureWarning: The default value of numeric_only in DataFrame.corr is deprecated. In a future vers
sns.heatmap(titanic_data.corr(), cmap="YlGnBu")



```

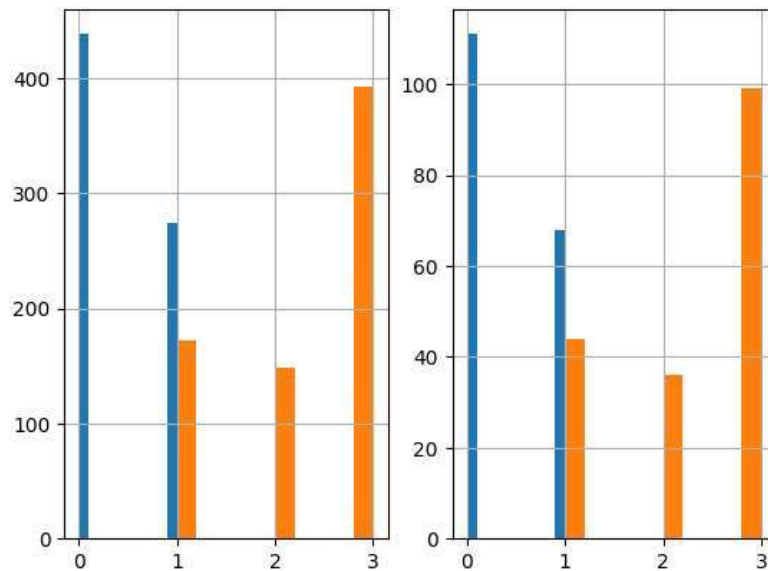
1 from sklearn.model_selection import StratifiedShuffleSplit
2
3 split=StratifiedShuffleSplit(n_splits=1,test_size=0.2)
4 for train_indices,test_indices in split.split(titanic_data,titanic_data[["Survived","Pclass","Sex"]]):
5     strat_train_set=titanic_data.loc[train_indices]
6     strat_test_set=titanic_data.loc[test_indices]

```

```

1 plt.subplot(1,2,1)
2 strat_train_set['Survived'].hist()
3 strat_train_set['Pclass'].hist()
4
5 plt.subplot(1,2,2)
6 strat_test_set['Survived'].hist()
7 strat_test_set['Pclass'].hist()
8
9 plt.show()

```



```
1 strat_train_set.info()
```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 712 entries, 871 to 459
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  -
0   PassengerId  712 non-null    int64
1   Survived     712 non-null    int64
2   Pclass       712 non-null    int64
3   Name         712 non-null    object
4   Sex          712 non-null    object
5   Age         575 non-null    float64
6   SibSp        712 non-null    int64
7   Parch        712 non-null    int64
8   Ticket       712 non-null    object
9   Fare         712 non-null    float64
10  Cabin        161 non-null    object
11  Embarked     710 non-null    object
dtypes: float64(2), int64(5), object(5)
memory usage: 72.3+ KB

```

Estimators

```

1 from sklearn.base import BaseEstimator, TransformerMixin
2 from sklearn.impute import SimpleImputer
3
4 class AgeImputer(BaseEstimator, TransformerMixin):
5
6     def fit(self, X, y=None):
7         return self
8
9     def transform(self, X):
10         imputer=SimpleImputer(strategy="mean")
11         X['Age']=imputer.fit_transform(X[['Age']])
12         return X

```

Encoding

```

1 from sklearn.preprocessing import OneHotEncoder
2
3 class FeatureEncoder(BaseEstimator, TransformerMixin):
4
5     def fit(self, X, y=None):
6         return self
7
8     def transform(self, X):
9         encoder=OneHotEncoder()
10        matrix=encoder.fit_transform(X[['Embarked']]).toarray()
11
12        column_names=["C","S","Q","N"]
13
14        for i in range(len(matrix.T)):
15            X[column_names[i]]=matrix.T[i]
16
17        matrix = encoder.fit_transform(X[['Sex']]).toarray()
18
19        column_names = ["Female", "Male"]
20
21        for i in range(len(matrix.T)):
22            X[column_names[i]] = matrix.T[i]
23
24        return X

```

Feature dropper

```

1 class FeatureDropper(BaseEstimator, TransformerMixin):
2
3     def fit(self, X, y=None):
4         return self
5
6     def transform(self, X):
7         return X.drop(["Embarked", "Name", "Ticket", "Cabin", "Sex","N"], axis=1, errors="ignore")
8

```

Pipeline

```

1 from sklearn.pipeline import Pipeline
2
3 pipeline = Pipeline([("ageimputer", AgeImputer()),
4                      ("featureencoder", FeatureEncoder()),
5                      ("featuredopper", FeatureDropper())])
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100

```

```

1 strat_train_set = pipeline.fit_transform(strat_train_set)

```

```

1 strat_train_set

```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	C	S	Q	Female	Male
871	872	1	1	47.000000	1	1	52.5542	0.0	0.0	1.0	1.0	0.0
400	401	1	3	39.000000	0	0	7.9250	0.0	0.0	1.0	0.0	1.0
786	787	1	3	18.000000	0	0	7.4958	0.0	0.0	1.0	1.0	0.0
6	7	0	1	54.000000	0	0	51.8625	0.0	0.0	1.0	0.0	1.0
127	128	1	3	24.000000	0	0	7.1417	0.0	0.0	1.0	0.0	1.0
...
879	880	1	1	56.000000	0	1	83.1583	1.0	0.0	0.0	1.0	0.0
487	488	0	1	58.000000	0	0	29.7000	1.0	0.0	0.0	0.0	1.0
123	124	1	2	32.500000	0	0	13.0000	0.0	0.0	1.0	1.0	0.0
522	523	0	3	30.095513	0	0	7.2250	1.0	0.0	0.0	0.0	1.0
459	460	0	3	30.095513	0	0	7.7500	0.0	1.0	0.0	0.0	1.0

712 rows × 12 columns

1 strat_train_set.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 712 entries, 871 to 459
Data columns (total 12 columns):
#   Column      Non-Null Count  Dtype
---  -
0   PassengerId  712 non-null    int64
1   Survived     712 non-null    int64
2   Pclass       712 non-null    int64
3   Age          712 non-null    float64
4   SibSp        712 non-null    int64
5   Parch        712 non-null    int64
6   Fare         712 non-null    float64
7   C            712 non-null    float64
8   S            712 non-null    float64
9   Q            712 non-null    float64
10  Female       712 non-null    float64
11  Male         712 non-null    float64
dtypes: float64(7), int64(5)
memory usage: 72.3 KB
```

```
1 from sklearn.preprocessing import StandardScaler
2
3 X = strat_train_set.drop(['Survived'], axis=1)
4 y = strat_train_set['Survived']
5
6 scaler = StandardScaler()
7 X_data = scaler.fit_transform(X)
8 y_data = y.to_numpy()
9
```

Model selection

```
1 from sklearn.ensemble import RandomForestClassifier
2 from sklearn.model_selection import GridSearchCV
3
4 clf = RandomForestClassifier()
5
6 param_gird = [
7     {"n_estimators" : [10, 100, 200, 500], "max_depth" : [None, 5, 10], "min_samples_split" : [2,3,4]}
8 ]
9
10 grid_search = GridSearchCV(clf, param_gird, cv=3, scoring="accuracy", return_train_score=True)
11 grid_search.fit(X_data, y_data)
12
```

```
GridSearchCV
GridSearchCV(cv=3, estimator=RandomForestClassifier(),
param_grid=[{'max_depth': [None, 5, 10],
'min_samples_split': [2, 3, 4],
'n_estimators': [10, 100, 200, 500]}],
return_train_score=True, scoring='accuracy')
  estimator: RandomForestClassifier
    RandomForestClassifier()
      RandomForestClassifier()
```

Result

```
1 final_clf = grid_search.best_estimator_
```

```
1 final_clf
```

```
RandomForestClassifier
RandomForestClassifier(max_depth=5, min_samples_split=4)
```

```
1 strat_test_set = pipeline.fit_transform(strat_test_set)
```

```
1 X_test = strat_test_set.drop(['Survived'], axis=1)
2 y_test = strat_test_set['Survived']
3
4 scaler = StandardScaler()
5 X_data_test = scaler.fit_transform(X_test)
6 y_data_test = y_test.to_numpy()
```

Model Accuracy

```
1 final_clf.score(X_data_test, y_data_test)
```

0.8324022346368715

```
1 final_data = pipeline.fit_transform(titanic_data)
```

```
1 final_data
```

	PassengerId	Survived	Pclass	Age	SibSp	Parch	Fare	C	S	Q	Female
0	1	0	3	22.000000	1	0	7.2500	0.0	0.0	1.0	0.0
1	2	1	1	38.000000	1	0	71.2833	1.0	0.0	0.0	1.0
2	3	1	3	26.000000	0	0	7.9250	0.0	0.0	1.0	1.0
3	4	1	1	35.000000	1	0	53.1000	0.0	0.0	1.0	1.0
4	5	0	3	35.000000	0	0	8.0500	0.0	0.0	1.0	0.0
...
886	887	0	2	27.000000	0	0	13.0000	0.0	0.0	1.0	0.0
887	888	1	1	19.000000	0	0	30.0000	0.0	0.0	1.0	1.0
888	889	0	3	29.699118	1	2	23.4500	0.0	0.0	1.0	1.0
889	890	1	1	26.000000	0	0	30.0000	1.0	0.0	0.0	0.0
890	891	0	3	32.000000	0	0	7.7500	0.0	1.0	0.0	0.0

891 rows × 12 columns



```
1 X_final = final_data.drop(['Survived'], axis=1)
2 y_final = final_data['Survived']

1 prod_clf = RandomForestClassifier()
2
3 param_gird = [
4     {"n_estimators" : [10, 100, 200, 500], "max_depth" : [None, 5, 10], "min_samples_split" : [2,3,4]}
5 ]
6
7 grid_search = GridSearchCV(prod_clf, param_gird, cv=3, scoring="accuracy", return_train_score=True)
8 grid_search.fit(X_data_final, y_data_final)
9
```

▼

GridSearchCV

GridSearchCV(cv=3, estimator=RandomForestClassifier(),
param_grid=[{'max_depth': [None, 5, 10],
'min_samples_split': [2, 3, 4],
'n_estimators': [10, 100, 200, 500]}],
return_train_score=True, scoring='accuracy')
▼ estimator: RandomForestClassifier
RandomForestClassifier()
▼ RandomForestClassifier
RandomForestClassifier()

Test data

```
1 prod_final_clf = grid_search.best_estimator_
```

Import Test dataset(csv file)

```
1 titanic_test_data = pd.read_csv("/content/test.csv")




1 final_test_data = pipeline.fit_transform(titanic_test_data)

1 X_final_test = final_test_data
2 X_final_test = X_final_test.fillna(method="ffill")
3
4 scaler = StandardScaler()
5 X_data_final_test = scaler.fit_transform(X_final_test)

1 predictions = prod_final_clf.predict(X_data_final_test)

1 final_df = pd.DataFrame(titanic_test_data['PassengerId'])
2 final_df['Survived'] = predictions

1 final_df
```

	PassengerId	Survived	
0	892	0	
1	893	0	
2	894	0	
3	895	0	
4	896	1	
...	
413	1305	0	
414	1306	1	
415	1307	0	