Machine Learning Lab - 2

DATA PREPROCESSING

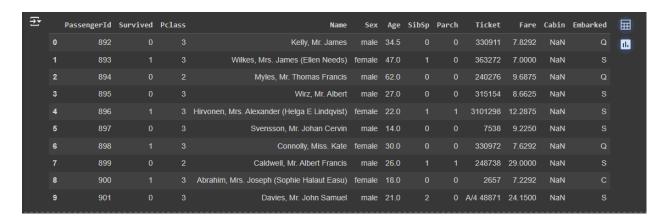
To apply various machine learning techniques to answer different questions, you should start by cleaning and preprocessing the data for each classification problem. For instance, with the Titanic dataset, you can prepare the data to be suitable for training and evaluating models for predicting or inferring demographic characteristics. Use this dataset to address eight distinct classification problems given below, ensuring that the data is properly preprocessed and that suitable features are selected for each model.

Instructions:

- 1. Data Preprocessing:
 - Load the Titanic dataset and clean basic data.
 - Handle missing values appropriately.
 - Convert categorical variables into numerical representations.

Output -

```
# Importing Dataset
import pandas as pd
data = pd.read_excel("/content/drive/MyDrive/Lab/ML/Lab_2/titanic.xlsx")
data.head(10)
```



```
missing_data = data.isna().sum()
missing_data
```

```
Passengerld
            0
 Survived
            0
  Pclass
            0
  Name
            0
  Sex
           0
           86
  Age
  SibSp
           0
  Parch
            0
  Ticket 0
  Fare
  Cabin 327
Embarked
            0
```

```
# Age Null Values Handling
mean_age = data["Age"].mean()
mean_age = int(mean_age)

data["Age"] = data["Age"].apply(lambda ele : mean_age if pd.isna(ele) else
ele)
data.head(50)
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	892			Kelly, Mr. James	male	34.5			330911	7.8292	NaN	Q
1	893			Wilkes, Mrs. James (Ellen Needs)	female	47.0			363272	7.0000	NaN	S
2	894			Myles, Mr. Thomas Francis	male	62.0			240276	9.6875	NaN	Q
3	895			Wirz, Mr. Albert	male	27.0			315154	8.6625	NaN	S
4	896			Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0			3101298	12.2875	NaN	S
5	897			Svensson, Mr. Johan Cervin	male	14.0			7538	9.2250	NaN	S
6	898			Connolly, Miss. Kate	female	30.0			330972	7.6292	NaN	Q
7	899			Caldwell, Mr. Albert Francis	male	26.0			248738	29.0000	NaN	s
8	900	1	3	Abrahim, Mrs. Joseph (Sophie Halaut Easu)	female	18.0	0	0	2657	7.2292	NaN	С

```
# Fare Null Values Handling
mean_fare = data["Fare"].mean()
```

```
mean_fare = int(mean_fare)

data["Fare"] = data["Fare"].apply(lambda ele : mean_fare if pd.isna(ele)

else ele)

data.head(20)
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
0	892			Kelly, Mr. James	male	34.5			330911	7.8292	NaN	Q
1	893			Wilkes, Mrs. James (Ellen Needs)	female	47.0			363272	7.0000	NaN	s
2	894			Myles, Mr. Thomas Francis	male	62.0			240276	9.6875	NaN	Q
3	895			Wirz, Mr. Albert	male	27.0			315154	8.6625	NaN	s
4	896			Hirvonen, Mrs. Alexander (Helga E Lindqvist)	female	22.0			3101298	12.2875	NaN	s
5	897			Svensson, Mr. Johan Cervin	male	14.0			7538	9.2250	NaN	s
6	898			Connolly, Miss. Kate	female	30.0			330972	7.6292	NaN	Q
7	899			Caldwell, Mr. Albert Francis	male	26.0			248738	29.0000	NaN	s
8	900			Abrahim, Mrs. Joseph (Sophie Halaut Easu)	female	18.0			2657	7.2292	NaN	С

```
# HasCabin Label Encoding
data["HasCabin"] = data["Cabin"].apply(lambda ele : 0 if pd.isna(ele) else
1)
```

```
# Sex Label Encoding
# male -> 0
# female -> 1
data["Sex"] = data["Sex"].map(lambda ele : 0 if ele == 'male' else 1)
```

```
# Embarked Label Encoding

# C -> 0

# Q -> 1

# s -> 2

data["Embarked"] = data["Embarked"].map({ 'C' : 0, 'Q' : 1, 'S' : 2 })
```

```
# Family Status Label Encoding
data["FamilyStatus"] = data.apply(lambda row : 1 if row["SibSp"] > 0 or
row["Parch"] > 0 else 0, axis = 1)
```

```
# HighFare Label Encoding
data["HighFare"] = data.apply(lambda row : 1 if row["Fare"] > mean_fare
else 0, axis = 1)
```

```
# Age Group Label Encoding
max_age = data["Age"].max() # 76.0
```

```
min_age = data["Age"].min() # 0.17

bins = [0, 10, 18, 40, 70, 100]

labels = [0, 1, 2, 3, 4]

data['AgeGroup'] = pd.cut(data['Age'], bins=bins, labels=labels)
```

```
# Name Label Encoding
uniqueNames = data["Name"].unique()
nameNumberMap = {}
for i, name in enumerate(uniqueNames):
   nameNumberMap[name] = i

data["Name"] = data["Name"].map(nameNumberMap)
```

```
# Ticket Label Encoding
tickets = data["Ticket"]
encodedTickets = []
for ticket in tickets:
  encodedTickets.append(str(ticket)[-5:])

data["Ticket"] = encodedTickets;
```

```
# Cabin Label Encoding
data["Cabin"] = data["Cabin"].map(lambda ele : "Unknown" if pd.isna(ele)
else ele)

data.head(15)

uniqueCabins = data["Cabin"].unique()
cabinNumberMap = {}
for i, name in enumerate(uniqueCabins):
    cabinNumberMap[name] = i

data["Cabin"] = data["Cabin"].map(cabinNumberMap)
```

	PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked	HasCabin	FamilyStatus	HighFare	AgeGroup
0	892					34.5			30911	7.8292						
1	893					47.0			63272	7.0000						
2	894					62.0			40276	9.6875						
3	895					27.0			15154	8.6625						
4	896					22.0			01298	12.2875						
5	897					14.0			7538	9.2250						
6	898					30.0			30972	7.6292						
7	899					26.0			48738	29.0000						
8	900					18.0			2657	7.2292						
9	901					21.0			48871	24.1500						