

```
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
import warnings
warnings.filterwarnings("ignore")

In [76]: hospital_data = pd.read_csv("hospital_appointment.csv")

In [77]: hospital_data
Out[77]:
   PatientID  AppointmentID  Gender  ScheduledDay  AppointmentDay  Age  Neighbourhood  Scholarship  Hypertension  Diabetes  Alcoholism  Handicap  SMS_received  No-show
0  2.96750e+13  5642903    F  2016-04-29T18:38:02Z  2016-04-29T00:00:00Z  62  JARDIM DA PENHA  0  0  1  0  0  0  0  0  No
1  5.59960e+14  5642903    M  2016-04-29T18:08:27Z  2016-04-29T00:00:00Z  56  JARDIM DA PENHA  0  0  0  0  0  0  0  0  No
2  4.32900e+12  5642949    F  2016-04-29T16:19:04Z  2016-04-29T00:00:00Z  62  MATA DA PRATA  0  0  0  0  0  0  0  0  No
3  8.67610e+12  5620828    F  2016-04-29T17:29:31Z  2016-04-29T00:00:00Z  8  PONTAL DE CAMBURU  0  0  0  0  0  0  0  0  No
4  8.84110e+12  5642949    F  2016-04-29T16:07:23Z  2016-04-29T00:00:00Z  56  JARDIM DA PENHA  0  0  1  1  0  0  0  0  No
...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...  ...
119522  2.57130e+12  5651768    F  2016-06-03T09:15:35Z  2016-06-07T00:00:00Z  56  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119523  3.596270e+12  5650093    F  2016-06-03T07:27:32Z  2016-06-07T00:00:00Z  51  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119524  1.56760e+13  5630692    F  2016-04-27T16:03:52Z  2016-06-07T00:00:00Z  21  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119525  9.213490e+13  5630323    F  2016-04-27T15:09:23Z  2016-06-07T00:00:00Z  38  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119526  3.775120e+14  5629448    F  2016-04-27T13:30:56Z  2016-06-07T00:00:00Z  54  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119527 rows x 14 columns

In [4]: hospital_data.tail()
Out[4]:
   PatientID  AppointmentID  Gender  ScheduledDay  AppointmentDay  Age  Neighbourhood  Scholarship  Hypertension  Diabetes  Alcoholism  Handicap  SMS_received  No-show
119522  2.57130e+12  5651768    F  2016-06-03T09:15:35Z  2016-06-07T00:00:00Z  56  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119523  3.596270e+12  5650093    F  2016-06-03T07:27:32Z  2016-06-07T00:00:00Z  51  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119524  1.56760e+13  5630692    F  2016-04-27T16:03:52Z  2016-06-07T00:00:00Z  21  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119525  9.213490e+13  5630323    F  2016-04-27T15:09:23Z  2016-06-07T00:00:00Z  38  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No
119526  3.775120e+14  5629448    F  2016-04-27T13:30:56Z  2016-06-07T00:00:00Z  54  MARIA ORTIZ  0  0  0  0  0  0  0  0  1  No

In [5]: hospital_data.shape
Out[5]:
(119527, 14)

In [6]: hospital_data.columns
Out[6]:
Index(['PatientID', 'AppointmentID', 'Gender', 'ScheduledDay', 'AppointmentDay', 'Age', 'Neighbourhood', 'Scholarship', 'Hypertension', 'Diabetes', 'Alcoholism', 'Handicap', 'SMS_received', 'No-show'],
      dtype='object')

In [53]: hospital_data.rename(columns={'No-show': 'No_show'}, inplace=True)
hospital_data.rename(columns={'Hypertension': 'Hypertension'}, inplace=True)
hospital_data.rename(columns={'Handicap': 'Handicap'}, inplace=True)

In [54]: hospital_data.info()
<class 'pandas.core.frame.DataFrame'>
Int64Index: 119527 entries, 0 to 119526
Data columns (total 14 columns):
#  Column                Non-Null Count  Dtype
---  ---
0  PatientID              119527 non-null  float64
1  AppointmentID          119527 non-null  int64
2  Gender                 119527 non-null  object
3  ScheduledDay           119527 non-null  object
4  AppointmentDay         119527 non-null  object
5  Age                    119527 non-null  int64
6  Neighbourhood          119527 non-null  object
7  Scholarship            119527 non-null  int64
8  Hypertension           119527 non-null  int64
9  Diabetes               119527 non-null  int64
10  Alcoholism             119527 non-null  int64
11  Handicap               119527 non-null  int64
12  SMS_received           119527 non-null  int64
13  No_show                119527 non-null  object
dtype: object

In [55]: hospital_data.describe()
Out[55]:
   PatientID  AppointmentID      Age  Scholarship  Hypertension  Diabetes  Alcoholism  Handicap  SMS_received
count  119527e+05  1.19527e+06  119527.000000  119527.000000  119527.000000  119527.000000  119527.000000  119527.000000
mean  1.474603e+14  5.678350e+06  37.088874  0.098266  0.197246  0.071985  0.030460  0.02248  0.321026
std  2.560949e+14  7.128575e+04  23.110205  0.297675  0.397921  0.256265  0.171666  0.161543  0.466873
min  3.920000e+14  5.030230e+06  -1.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000
25%  4.172615e+12  5.640286e+06  18.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000
50%  3.173280e+13  5.680573e+06  37.000000  0.000000  0.000000  0.000000  0.000000  0.000000  0.000000
75%  4.849170e+13  5.725244e+06  55.000000  0.000000  0.000000  0.000000  0.000000  0.000000  1.000000
max  8.999829e+14  5.749049e+06  115.000000  1.000000  1.000000  1.000000  1.000000  4.000000  1.000000

In [56]: hospital_data.isnull().sum()
Out[56]:
PatientID      0
AppointmentID  0
Gender          0
ScheduledDay    0
AppointmentDay  0
Age             0
Neighbourhood  0
Scholarship     0
Hypertension    0
Diabetes        0
Alcoholism      0
Handicap        0
SMS_received    0
No_show         0
dtype: int64

In [57]: hospital_data.nunique()
Out[57]:
PatientID      41344
AppointmentID   11827
Gender          27
ScheduledDay   18349
AppointmentDay  164
Age            164
Neighbourhood  61
Scholarship    2
Hypertension   2
Diabetes       2
Alcoholism     2
Handicap       6
SMS_received   2
No_show        2
dtype: int64

In [58]: hospital_data1 = hospital_data[['gender', 'Scholarship', 'Hypertension', 'Diabetes', 'Alcoholism', 'Handicap', 'SMS_received', 'No_show']]

In [59]: for i in hospital_data1.columns:
plt.figure(figsize=(15,6))
sns.countplot(hospital_data1[i], data = hospital_data1, palette='hls')
plt.show()

In [60]: hospital_data2 = hospital_data.copy()
In [61]: hospital_data2.drop(['PatientID', 'AppointmentID', 'ScheduledDay', 'AppointmentDay'], axis=1)
Out[61]:
   Gender  Age  Neighbourhood  Scholarship  Hypertension  Diabetes  Alcoholism  Handicap  SMS_received  No-show
0  F  62  JARDIM DA PENHA  0  0  1  0  0  0  0  No
1  M  56  JARDIM DA PENHA  0  0  0  0  0  0  0  No
2  F  62  MATA DA PRATA  0  0  0  0  0  0  0  No
3  F  8  PONTAL DE CAMBURU  0  0  0  0  0  0  0  No
4  F  56  JARDIM DA PENHA  0  1  1  0  0  0  0  No
...  ...  ...  ...  ...  ...  ...  ...  ...  ...
119522  F  51  MARIA ORTIZ  0  0  0  0  0  0  0  1  No
119523  F  51  MARIA ORTIZ  0  0  0  0  0  0  0  1  No
119524  F  21  MARIA ORTIZ  0  0  0  0  0  0  0  1  No
119525  F  38  MARIA ORTIZ  0  0  0  0  0  0  0  1  No
119526  F  54  MARIA ORTIZ  0  0  0  0  0  0  0  1  No
119527 rows x 10 columns

In [62]: hospital_data2['Age'] = replace(hospital_data2['Age'].mean(), inplace = True)

In [18]: hospital_data2['Age'] = hospital_data2['Age'].abs()

In [19]: hospital_data2.hist(figsize=(12,12))

In [20]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'Gender', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [21]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'Diabetes', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [22]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'Alcoholism', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [23]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'Hypertension', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [24]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'Alcoholism', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [25]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'Alcoholism', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [26]: plt.figure(figsize=(15,6))
sns.countplot('No_show', hue = 'SMS_received', data = hospital_data2, palette='hls')
plt.xticks(rotation = 90)
plt.show()

In [27]: plt.figure(figsize=(20,20))
e = hospital_data2.groupby(['Neighbourhood', 'No_show']).size().unstack()
e.ys.plot(kind='bar', alpha=.5, color='red', label='no show')
e.xs.plot(kind='bar', alpha=.5, color='green', label='show')
plt.legend()
plt.title("The relation between neighbourhood and showing up")
plt.xlabel("Neighbourhood")
plt.ylabel("patients")

In [28]: plt.figure(figsize=(15,6))
sns.histplot(hospital_data2[['Age']])
plt.show()

In [29]: hospital_age = hospital_data2['Age']
Q3 = hospital_age.quantile(0.75)
Q1 = hospital_age.quantile(0.25)
IQR = Q3-Q1
lower_limit = Q1 - (1.5*IQR)
upper_limit = Q3 + (1.5*IQR)
age_outliers = hospital_age[(hospital_age < lower_limit) | (hospital_age > upper_limit)]

In [29]: 63912  115.8
63913  115.8
68237  115.8
76284  115.8
97686  115.8
Name: Age, dtype: float64

In [30]: lower_limit
Out[30]:
-32.5

In [31]: upper_limit
Out[31]:
107.5

In [30]: hospital_data_new = hospital_data.drop([63912, 63913, 68127, 76284, 97686])

In [31]: hospital_data_new.head()
Out[31]:
   PatientID  AppointmentID      Age  Scholarship  Hypertension  Diabetes  Alcoholism  Handicap  SMS_received
0  2.96750e+13  5642903    F  2016-04-29T18:38:02Z  2016-04-29T00:00:00Z  62  JARDIM DA PENHA  0  0  1  0  0  0  0  No
1  5.59960e+14  5642903    M  2016-04-29T18:08:27Z  2016-04-29T00:00:00Z  56  JARDIM DA PENHA  0  0  0  0  0  0  0  No
2  4.32900e+12  5642949    F  2016-04-29T16:19:04Z  2016-04-29T00:00:00Z  62  MATA DA PRATA  0  0  0  0  0  0  0  No
3  8.67610e+12  5620828    F  2016-04-29T17:29:31Z  2016-04-29T00:00:00Z  8  PONTAL DE CAMBURU  0  0  0  0  0  0  0  No
4  8.84110e+12  5642949    F  2016-04-29T16:07:23Z  2016-04-29T00:00:00Z  56  JARDIM DA PENHA  0  0  1  1  0  0  0  No

In [32]: from sklearn import preprocessing
label_encoder = preprocessing.LabelEncoder()
hospital_data_new['Gender'] = label_encoder.fit_transform(hospital_data_new['Gender'])
hospital_data_new['No_show'] = label_encoder.fit_transform(hospital_data_new['No_show'])

In [33]: hospital_data_new['No_show']
Out[33]:
0  0
1  0
2  0
3  0
4  0
...
119522  0
119523  0
119524  0
119525  0
119526  0
Name: No_show, Length: 119522, dtype: int64

In [35]: hospital_data_new.columns
Out[35]:
Index(['PatientID', 'AppointmentID', 'Gender', 'ScheduledDay', 'AppointmentDay', 'Age', 'Neighbourhood', 'Scholarship', 'Hypertension', 'Diabetes', 'Alcoholism', 'Handicap', 'SMS_received', 'No_show'],
      dtype='object')

In [36]: x = hospital_data_new[['Gender', 'Scholarship', 'Hypertension', 'Diabetes', 'Alcoholism', 'Handicap', 'SMS_received']]

In [37]: x = hospital_data_new.No_show

In [38]: x.shape
Out[38]:
(119522, 7)

In [39]: y.shape
Out[39]:
(119522,)

In [40]: from sklearn.model_selection import train_test_split

In [41]: x_train, x_test, y_train, y_test = train_test_split(x, y, test_size=0.2, random_state=42)

In [42]: from sklearn.linear_model import LogisticRegression

In [43]: classifier = LogisticRegression(random_state=0)
classifier.fit(x_train, y_train)

In [44]: print("Training Accuracy :", classifier.score(x_train, y_train))
print("Testing Accuracy :", classifier.score(x_test, y_test))
Training Accuracy : 0.798364547899392
Testing Accuracy : 0.749699113358516

In [45]: from sklearn.tree import DecisionTreeClassifier

In [46]: classifier_dt = DecisionTreeClassifier(criterion='entropy', random_state=0)
classifier_dt.fit(x_train, y_train)

In [47]: print("Training Accuracy :", classifier_dt.score(x_train, y_train))
print("Testing Accuracy :", classifier_dt.score(x_test, y_test))
Training Accuracy : 0.788602878785738
Testing Accuracy : 0.7779504888683381

In [ ]:
```