Pclass 0 Name 0 0 Sex Age 177 0 SibSp Parch 0 0 Ticket 0 Fare 687 Cabin Embarked 2 dtype: int64 The columns having null values are: Age, Cabin, Embarked. They need to be filled up with appropriate values later on. Features: The titanic dataset has roughly the following types of features: Categorical/Nominal: Variables that can be divided into multiple categories but having no order or priority. Eg. Embarked (C = Cherbourg; Q = Queenstown; S = Southampton) Binary: A subtype of categorical features, where the variable has only two categories. Eg: Sex (Male/Female) Ordinal: They are similar to categorical features but they have an order (i.e can be sorted). Eg. Pclass (1, 2, 3) Continuous: They can take up any value between the minimum and maximum values in a column. Eg. Age, Fare Count: They represent the count of a variable. Eg. SibSp, Parch Useless: They don't contribute to the final outcome of an ML model. Here, Passengerld, Name, Cabin and Ticket might fall into this category. **Graphical Analysis** In [5]: import seaborn as sns import matplotlib.pyplot as plt sns.catplot(x ="Sex", hue ="Survived", kind ="count", data = titanic) <seaborn.axisgrid.FacetGrid at 0x2c39e7aea00> 300

Just by observing the graph, it can be approximated that the survival rate of men is around 20% and that of women is around 75%. Therefore, whether a passenger is a male or a female plays an

It helps in determining if higher-class passengers had more survival rate than the lower class ones or vice versa. Class 1 passengers have a higher survival chance compared to classes 2 and 3. It

D:\ANACONDA\lib\site-packages\seaborn\categorical.py:3717: UserWarning: The `factorplot` function has been renamed to `catplot`. The original name will be rem

100 -

male

0

In [6]:

Out[6]:

In [7]:

In [8]:

Out[8]:

In [9]:

Out[9]:

Divide Fare into 4 bins

on the height of bars.

data = titanic)

0.2

0.1

250

T 150

100

50

In []:

THANK YOU

titanic['Fare_Range'] = pd.qcut(titanic['Fare'], 4)

<AxesSubplot:xlabel='Fare_Range', ylabel='Survived'>

(-0.001, 7.91] (7.91, 14.454] (14.454, 31.0] (31.0, 512.329] Fare_Range

passenger paid a higher fare, the survival rate is more.

Barplot - Shows approximate values based

sns.barplot(x ='Fare_Range', y ='Survived',

0.8

Loading data using Pandas

3

3

1

3

titanic = pd.read_csv('https://raw.githubusercontent.com/dsrscientist/dataset1/master/titanic_train.csv')

Braund, Mr. Owen Harris

Futrelle, Mrs. Jacques Heath (Lily May Peel) female 35.0

Allen, Mr. William Henry

1 Cumings, Mrs. John Bradley (Florence Briggs Th... female 38.0

Survived

Name

Heikkinen, Miss. Laina female

Sex Age SibSp

22.0

26.0

male 35.0

male

Parch

0

1

0

Ticket

PC 17599 71.2833

113803 53.1000

A/5 21171

373450

0 STON/O2. 3101282

Fare

7.2500

7.9250

8.0500

Cabin Embarked

S

С

S

S

S

NaN

C85

NaN

C123

NaN

import pandas as pd

1

3

titanic.isnull().sum()

PassengerId

Survived

Passengerld Survived Pclass

1

0

0

0

Checking the NULL values

titanic.head()

In [3]:

Out[3]:

Out[4]:

0

2

- 350 - 300 - 300 - 250 - 200 - 150 - 100

1

Age (Continuous Feature) vs Survived

Survived

Violinplot Displays distribution of data

across all levels of a category.

data = titanic, split = True)

Adding a column Family_Size
titanic['Family_Size'] = 0

Factorplot for Family_Size

Factorplot for Alone

implies that Pclass contributes a lot to a passenger's survival rate.

sns.violinplot(x ="Sex", y ="Age", hue ="Survived",

female

Sex

group = titanic.groupby(['Pclass', 'Survived'])

<AxesSubplot:xlabel='Survived', ylabel='Pclass'>

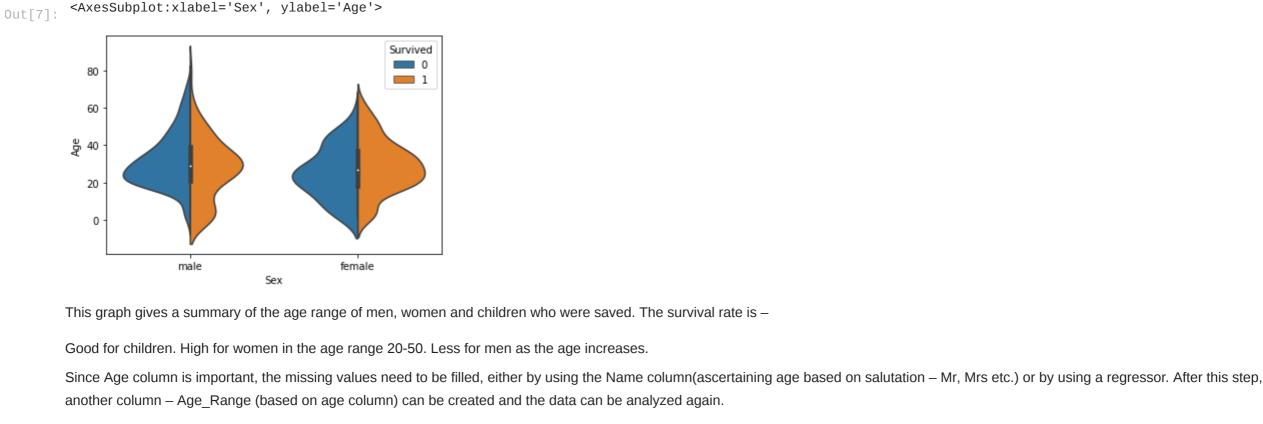
Heatmap - Color encoded 2D representation of data.
sns.heatmap(pclass_survived, annot = True, fmt = "d")

Pclass (Ordinal Feature) vs Survived

Group the dataset by Pclass and Survived and then unstack them

important role in determining if one is going to survive.

pclass_survived = group.size().unstack()



Adding a column Alone
titanic['Alone'] = 0
titanic.loc[titanic.Family_Size == 0, 'Alone'] = 1

titanic['Family_Size'] = titanic['Parch']+titanic['SibSp']

sns.factorplot(x ='Family_Size', y ='Survived', data = titanic)

sns.factorplot(x ='Alone', y ='Survived', data = titanic)

<seaborn.axisgrid.FacetGrid at 0x2c3a3b5c400>

Factor plot for Family_Size (Count Feature) and Family Size.

oved in a future release. Please update your code. Note that the default `kind` in `factorplot` (`'point'`) has changed `'strip'` in `catplot`.
warnings.warn(msg)
D:\ANACONDA\lib\site-packages\seaborn\categorical.py:3717: UserWarning: The `factorplot` function has been renamed to `catplot`. The original name will be rem
oved in a future release. Please update your code. Note that the default `kind` in `factorplot` (`'point'`) has changed `'strip'` in `catplot`.
warnings.warn(msg)

Family_Size denotes the number of people in a passenger's family. It is calculated by summing the SibSp and Parch columns of a respective passenger. Also, another column Alone is added to check the chances of survival of a lone passenger against the one with a family.

Important observations—
If a passenger is alone, the survival rate is less. If the family size is greater than 5, chances of survival decrease considerably.

Bar Plot for Fare (Continuous Feature)

0.6 -0.5 -

Categorical Count Plots for Embarked Feature

In [10]: # Countplot sns.catplot(x = 'Embarked', hue = 'Survived', kind = 'count', col = 'Pclass', data = titanic)

Out[10]: <seaborn.axisgrid.FacetGrid at 0x2c3a3834e20>

Pclass = 1

Pclass = 2

Q

Some notable observations are:

Majority of the passengers boarded from S. So, the missing values can be filled with S. Majority of class 3 passengers boarded from Q. S looks lucky for class 1 and 2 passengers compared to class 3.

Conclusion:

The columns that can be dropped are:

PassengerId, Name, Ticket, Cabin: They are strings, cannot be categorized and don't contribute much to the outcome.

Age, Fare: Instead, the respective range columns are retained.

The titanic data can be analyzed using many more graph techniques and also more column correlations, than, as described in this article.

Embarked

Q

Fare denotes the fare paid by a passenger. As the values in this column are continuous, they need to be put in separate bins(as done for Age feature) to get a clear idea. It can be concluded that if a

Pclass = 3

Embarked

Survived

Once the EDA is completed, the resultant dataset can be used for predictions.

Embarked