

```
In [1]: import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
import seaborn as sns

In [2]: import matplotlib.pyplot as plt

In [3]: dataframe_train = pd.read_csv('titanic.csv')
dataframe_train.head()

Out[3]:
```

	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	Cummings, Mrs. John Bradley (Florence Briggs Th...	female	38.0	1	0	PC 17599	71.2833	C85	C
2	3	1	3	Heikinen, 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

```
In [4]: dataframe_train.shape

Out[4]: (891, 12)
```

```
In [5]: dataframe_train.dtypes

Out[5]: PassengerId    int64
Survived         int64
Pclass          object
Name            object
Sex             object
Age            float64
SibSp           int64
Parch           int64
Ticket          object
Fare           float64
Cabin          object
Embarked        object
dtype: object

In [6]: dataframe_train.Survived.value_counts()

Out[6]: 0    549
        1    342
        Name: Survived, dtype: int64

In [7]: dataframe_train.Sex.value_counts()

Out[7]: male      577
        female  314
        Name: Sex, dtype: int64

In [8]: dataframe_train.Embarked.value_counts()

Out[8]: S    644
        C    168
        Q      0
        Name: Embarked, dtype: int64

In [9]: dataframe_train.isnull().sum()

Out[9]: PassengerId    0
Survived          0
Pclass            0
Name              0
Age              177
SibSp             0
Parch            0
Ticket            0
Fare              0
Cabin           687
Embarked         2
dtype: int64

In [10]: dataframe_train.columns

Out[10]: Index(['PassengerId', 'Survived', 'Pclass', 'Name', 'Sex', 'Age', 'SibSp',
        'Parch', 'Ticket', 'Fare', 'Cabin', 'Embarked'],
        dtype='object')

In [11]: # All these columns seems to be irrelevant.
columns_to_drop = ['PassengerId', 'Name', 'Ticket', 'Cabin', 'Embarked']

In [12]: dataframe_train = dataframe_train.drop(columns_to_drop, axis = 1)
dataframe_train.head()

Out[12]:
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare
0	0	3	male	22.0	1	0	7.2500
1	1	1	female	38.0	1	0	71.2833
2	1	3	female	26.0	0	0	7.9250
3	1	1	female	35.0	1	0	53.1000
4	0	3	male	35.0	0	0	8.0500

```
In [13]: dataframe_train[dataframe_train['Age'].isna()]

Out[13]:
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare
5	0	3	male	NaN	0	0	8.4583
17	1	2	male	NaN	0	0	13.0000
19	1	3	female	NaN	0	0	7.2250
26	0	3	male	NaN	0	0	7.2250
28	1	3	female	NaN	0	0	7.8792
...	...	...	...	...	...	...	...
859	0	3	male	NaN	0	0	7.2292
863	0	3	female	NaN	8	2	69.5500
868	0	3	male	NaN	0	0	9.5000
878	0	3	male	NaN	0	0	7.8958
888	0	3	female	NaN	1	2	23.4500

```
177 rows x 7 columns

In [14]: # Smart way of replacing null values in age column by looking into Pclass as reference column
dataframe_train.groupby('Pclass').mean()[['Age']]

Out[14]:
```

Pclass	Age
1	30.230441
2	29.877830
3	25.140620

```
In [15]: dataframe_train.head()

Out[15]:
```

	Survived	Pclass	Sex	Age	SibSp	Parch	Fare
0	0	3	male	22.0	1	0	7.2500
1	1	1	female	38.0	1	0	71.2833
2	1	3	female	26.0	0	0	7.9250
3	1	1	female	35.0	1	0	53.1000
4	0	3	male	35.0	0	0	8.0500

```
In [16]: def age_approx(cols):
    Age = cols[0]
    Pclass = cols[1]
    if pd.isnull(Age):
        if Pclass == 1:
            return 39
        elif Pclass == 2:
            return 30
        else:
            return 24
    else:
        return Age

In [17]: dataframe_train['Age'] = dataframe_train[['Age', 'Pclass']].apply(age_approx, axis = 1)

In [18]: dataframe_train.isna().sum()

Out[18]: Survived      0
Pclass           0
Sex              0
Age              0
SibSp           0
Parch           0
Fare            0
dtype: int64

In [19]: dataframe_train.dtypes

Out[19]: Survived    int64
Pclass         int64
Sex            object
Age           float64
SibSp         int64
Parch         int64
Fare         float64
dtype: object

In [20]: # applying one hot encoding on sex col
dataframe_train_one_hot = pd.get_dummies(dataframe_train, columns = ['Sex'])

In [21]: dataframe_train_one_hot

Out[21]:
```

	Survived	Pclass	Age	SibSp	Parch	Fare	Sex_female	Sex_male
0	0	3	22.0	1	0	7.2500	0	1
1	1	1	38.0	1	0	71.2833	1	0
2	1	3	26.0	0	0	7.9250	1	0
3	1	1	35.0	1	0	53.1000	1	0
4	0	3	35.0	0	0	8.0500	0	1
...	...	...	...	...	...	...	...	...
886	0	2	27.0	0	0	13.0000	0	1
887	1	1	19.0	0	0	30.0000	1	0
888	0	3	24.0	1	2	23.4500	1	0
889	1	1	26.0	0	0	30.0000	0	1
890	0	3	32.0	0	0	7.7500	0	1

```
891 rows x 9 columns

In [22]: from sklearn.preprocessing import StandardScaler, MinMaxScaler
age_std_scaler = StandardScaler()
fare_std_scaler = StandardScaler()

dataframe_train_one_hot['Age'] = age_std_scaler.fit_transform(dataframe_train_one_hot[['Age']])
dataframe_train_one_hot['Fare'] = fare_std_scaler.fit_transform(dataframe_train_one_hot[['Fare']])

In [23]: dataframe_train_one_hot.head()

Out[23]:
```

	Survived	Pclass	Age	SibSp	Parch	Fare	Sex_female	Sex_male
0	0	3	-0.538003	1	0	-0.502445	0	1
1	1	1	0.666578	1	0	0.786845	1	0
2	1	3	-0.238559	0	0	-0.488854	1	0
3	1	1	0.440719	1	0	0.420730	1	0
4	0	3	0.440719	0	0	-0.486337	0	1

```
In [24]: dataframe_train_one_hot.columns

Out[24]: Index(['Survived', 'Pclass', 'Age', 'SibSp', 'Parch', 'Fare', 'Sex_female',
        'Sex_male'],
        dtype='object')

In [25]: #to find the pairwise correlation
plt.figure(figsize=(20,9));
sns.heatmap(dataframe_train_one_hot.corr(),annot=True,cmap='bwr');
plt.title("Correlation map")

Out[25]: Text(0.5, 1.0, 'Correlation map')
```

```
In [26]: X = dataframe_train_one_hot[['Pclass', 'Age', 'SibSp', 'Parch', 'Fare', 'Sex_female', 'Sex_male']]
y = dataframe_train_one_hot['Survived']

In [27]: X.shape, y.shape

Out[27]: ((891, 7), (891,))

In [28]: X.head()

Out[28]:
```

	Pclass	Age	SibSp	Parch	Fare	Sex_female	Sex_male
0	3	-0.538003	1	0	-0.502445	0	1
1	1	0.666578	1	0	0.786845	1	0
2	3	-0.238558	0	0	-0.488854	1	0
3	1	0.440719	1	0	0.420730	1	0
4	3	0.440719	0	0	-0.486337	0	1

```
In [29]: from sklearn.model_selection import train_test_split
X_train, X_test, y_train, y_test = train_test_split(X, y, random_state = 0)

In [30]: X_train.shape, X_test.shape, y_train.shape, y_test.shape

Out[30]: ((668, 7), (223, 7), (668, ), (223, ))

In [31]: from sklearn.linear_model import LogisticRegression
logreg = LogisticRegression()

In [32]: logreg.fit(X_train, y_train)

Out[32]: LogisticRegression()

In [33]: # X_test['Fare'] = fare_std_scaler.transform(X_test[['Fare']])
# X_test['Age'] = age_std_scaler.transform(X_test[['Age']])

In [34]: y_pred = logreg.predict(X_test)
y_pred

Out[34]: array([1, 0, 0, 1, 0, 0, 0, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 1, 0, 0,
        1, 0, 0, 0, 0, 0, 1, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 0, 0, 0,
        1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 1, 1, 0, 0,
        1, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 1, 0, 1, 0, 0, 1, 1, 0, 0, 0, 1,
        1, 1, 0, 0, 1, 1, 0, 0, 0, 0, 1, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
        0, 0, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
        0, 1, 0, 1, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1,
        0, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0, 0, 0, 1, 0,
        0, 0, 0, 0, 0, 0])

In [35]: for i in range(len(X_test)):
    print(logreg.predict_proba(X_test)[1])

Out[35]: (0.06944503 0.93055497)
(0.53465755 0.46534245)
(0.89714331 0.10285669)
(0.13578224 0.86421776)
(0.55866000 0.44133991)
(0.878918 0.121082)
(0.9183442 0.08165558)
(0.88031512 0.11964486)
(0.87285309 0.12794691)
(0.92295095 0.07704995)
(0.64421115 0.35571889)
(0.28386448 0.71613552)
(0.74749826 0.25250174)
(0.88015984 0.11964486)
(0.13234865 0.86765135)
(0.89323168 0.10668332)
(0.87282382 0.12791686)
(0.93918629 0.06089371)
(0.20599858 0.79409132)
(0.4217501 0.57824991)
(0.87285309 0.12794691)
(0.91368216 0.08631724)
(0.69877101 0.30122899)
(0.53799746 0.46200254)
(0.67184351 0.32815615)
(0.90713823 0.09286177)
(0.93629804 0.06376096)
(0.74749826 0.25250174)
(0.78554114 0.21445896)
(0.38284333 0.61795667)
(0.12618751 0.87381249)
(0.87992379 0.12067621)
(0.92911187 0.07088813)
(0.94351134 0.05634866)
(0.44874492 0.55125508)
(0.87221798 0.12782826)
(0.90531519 0.09454681)
(0.79125591 0.20874409)
(0.93469797 0.06530263)
(0.43411643 0.56508357)
(0.57963389 0.42038611)
(0.85337789 0.14062211)
(0.84717173 0.15028927)
(0.90393546 0.09686454)
(0.90882568 0.09197432)
(0.88538701 0.11461299)
(0.87284793 0.12795207)
(0.8847034 0.11552966)
(0.186885 0.813115)
(0.78514978 0.21485822)
(0.88518492 0.11403198)
(0.87285309 0.12794691)
(0.93949836 0.06095164)
(0.87625296 0.12347036)
(0.90450184 0.09549896)
(0.84184997 0.15815983)
(0.14752368 0.85247132)
(0.82916558 0.17083642)
(0.38626575 0.61373425)
(0.17471269 0.82520731)
(0.87185446 0.12814554)
(0.87140869 0.12859131)
(0.93738512 0.06261408)
(0.87282382 0.12797698)
(0.690303 0.309697)
(0.18688532 0.813115)
(0.77291527 0.22788473)
(0.43267815 0.56733985)
(0.73079605 0.26926935)
(0.44978169 0.55021831)
(0.87216845 0.12783955)
(0.65082827 0.34891713)
(0.75209637 0.24799363)
(0.86787136 0.13292864)
(0.80533332 0.19666666)
(0.25786211 0.74213789)
(0.19538856 0.80461364)
(0.87786993 0.12213007)
(0.89332727 0.10667273)
(0.86159991 0.13834099)
(0.90777776 0.09222224)
(0.36492866 0.63507134)
(0.86271088 0.13728912)
(0.6346581 0.3653419)
(0.3334547 0.6665453)
(0.8989811 0.10109184)
(0.8863265 0.11387329)
(0.36492331 0.63507669)
(0.38414712 0.61585288)
(0.13248481 0.86759152)
(0.69362811 0.30637189)
(0.87285309 0.12794691)
(0.68931428 0.31068572)
(0.1926451 0.8073549)
(0.5075458 0.4924542)
(0.85762859 0.14237941)
(0.57298046 0.42709116)
(0.87199121 0.12068079)
(0.2132274 0.7867726)
(0.87463462 0.12505158)
(0.85741975 0.14258925)
(0.49782713 0.50217253)
(0.87219738 0.12782626)
(0.77034373 0.22955627)
(0.49782713 0.50217253)
(0.93218741 0.06781259)
(0.89274663 0.10725337)
(0.85244168 0.14755832)
(0.4645479 0.5354521)
(0.22737722 0.7726276)
(0.94440636 0.05559364)
(0.94146896 0.05853184)
(0.95211458 0.04788542)
(0.90984896 0.09015104)
(0.42920373 0.57078627)
(0.73079605 0.26926935)
(0.84325983 0.15674017)
(0.87216845 0.12783955)
(0.89828272 0.10091126)
(0.3462568 0.6537432)
(0.74053187 0.25948813)
(0.35160653 0.64831547)
(0.86734292 0.13265768)
(0.75795937 0.24234063)
(0.88910441 0.11080959)
(0.20901934 0.79090866)
(0.89672421 0.90327379)
(0.87655882 0.12344918)
(0.91742657 0.08257343)
(0.88929221 0.11089779)
(0.90450184 0.09549896)
(0.84184997 0.15815983)
(0.14752368 0.85247132)
(0.82916558 0.17083642)
(0.38626575 0.61373425)
(0.17471269 0.82520731)
(0.87185446 0.12814554)
(0.87140869 0.12859131)
(0.93738512 0.06261408)
(0.87282382 0.12797698)
(0.690303 0.309697)
(0.18688532 0.813115)
(0.77291527 0.22788473)
(0.43267815 0.56733985)
(0.73079605 0.26926935)
(0.44978169 0.55021831)
(0.87216845 0.12783955)
(0.65082827 0.34891713)
(0.75209637 0.24799363)
(0.86787136 0.13292864)
(0.80533332 0.19666666)
(0.25786211 0.74213789)
(0.19538856 0.80461364)
(0.87786993 0.12213007)
(0.89332727 0.10667273)
(0.86159991 0.13834099)
(0.90777776 0.09222224)
(0.36492866 0.63507134)
(0.86271088 0.13728912)
(0.6346581 0.3653419)
(0.3334547 0.6665453)
(0.8989811 0.10109184)
(0.8863265 0.11387329)
(0.36492331 0.63507669)
(0.38414712 0.61585288)
(0.13248481 0.86759152)
(0.69362811 0.30637189)
(0.87285309 0.12794691)
(0.68931428 0.31068572)
(0.1926451 0.8073549)
(0.5075458 0.4924542)
(0.85762859 0.14237941)
(0.57298046 0.42709116)
(0.87199121 0.12068079)
(0.2132274 0.7867726)
(0.87463462 0.12505158)
(0.85741975 0.14258925)
(0.49782713 0.50217253)
(0.87219738 0.12782626)
(0.77034373 0.22955627)
(0.49782713 0.50217253)
(0.93218741 0.06781259)
(0.89274663 0.10725337)
(0.85244168 0.14755832)
(0.4645479 0.5354521)
(0.22737722 0.7726276)
(0.94440636 0.05559364)
(0.94146896 0.05853184)
(0.95211458 0.04788542)
(0.90984896 0.09015104)
(0.42920373 0.57078627)
(0.73079605 0.26926935)
(0.84325983 0.15674017)
(0.87216845 0.12783955)
(0.89828272 0.10091126)
(0.3462568 0.6537432)
(0.74053187 0.25948813)
(0.35160653 0.64831547)
(0.86734292 0.13265768)
(0.75795937 0.24234063)
(0.88910441 0.11080959)
(0.20901934 0.79090866)
(0.89672421 0.90327379)
(0.87655882 0.12344918)
(0.91742657 0.08257343)
(0.88929221 0.11089779)
(0.90450184 0.09549896)
(0.84184997 0.15815983)
(0.14752368 0.85247132)
(0.82916558 0.17083642)
(0.38626575 0.61373425)
(0.17471269 0.82520731)
(0.87185446 0.12814554)
(0.87140869 0.12859131)
(0.93738512 0.06261408)
(0.87282382 0.12797698)
(0.690303 0.309697)
(0.18688532 0.813115)
(0.77291527 0.22788473)
(0.43267815 0.56733985)
(0.73079605 0.26926935)
(0.44978169 0.55021831)
(0.87216845 0.12783955)
(0.65082827 0.34891713)
(0.75209637 0.24799363)
(0.86787136 0.13292864)
(0.80533332 0.19666666)
(0.25786211 0.74213789)
(0.19538856 0.80461364)
(0.87786993 0.12213007)
(0.89332727 0.10667273)
(0.86159991 0.13834099)
(0.90777776 0.09222224)
(0.36492866 0.63507134)
(0.86271088 0.13728912)
(0.6346581 0.3653419)
(0.3334547 0.6665453)
(0.8989811 0.10109184)
(0.8863265 0.11387329)
(0.36492331 0.63507669)
(0.38414712 0.61585288)
(0.13248481 0.86759152)
(0.69362811 0.30637189)
(0.87285309 0.12794691)
(0.68931428 0.31068572)
(0.1926451 0.8073549)
(0.5075458 0.4924542)
(0.85762859 0.14237941)
(0.57298046 0.42709116)
(0.87199121 0.12068079)
(0.2132274 0.7867726)
(0.87463462 0.12505158)
(0.85741975 0.14258925)
(0.49782713 0.50217253)
(0.87219738 0.12782626)
(0.77034373 0.22955627)
(0.49782713 0.50217253)
(0.93218741 0.06781259)
(0.89274663 0.10725337)
(0.8
```