

## Contents

<b>1</b>	<b>Introduction</b>	<b>1</b>
<b>2</b>	<b>Preparation</b>	<b>1</b>
<b>3</b>	<b>A first look at the dataset</b>	<b>3</b>
<b>4</b>	<b>First figures</b>	<b>5</b>
<b>5</b>	<b>Graphic illustration of a prediction</b>	<b>11</b>

In this paper we look at the chances of surviving the Titanic using python. We predict the chances of an unseen set of data by using supervised machinelearning on the known dataset.

We use information from [this site](#).

This is one of the first drafts to get to know the dataset and to experiment with the python and all the packages included.

## 1 Introduction

## 2 Preparation

---

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

---

The adventure begins with importing the right packages. The dataset is downloaded from [kaggle site](#) as csv\\_file. Next the data is read into a dataframe by using pandas' pd.read\\_csv

---

```
1 data = pd.read_csv('titanic.csv')
```

---

```
FileNotFoundErrorTraceback (most recent call last)
```

```
<ipython-input-2-74241d1367d4> in <module>()
```

```
----> 1 data = pd.read_csv('titanic.csv')
```

```
~/anaconda3/lib/python3.6/site-packages/pandas/io/parsers.py in parser_f(filepath_or_
707                                     skip_blank_lines=skip_blank_lines)
```

```

708
--> 709         return _read(filepath_or_buffer, kwds)
710
711     parser_f.__name__ = name

~/anaconda3/lib/python3.6/site-packages/pandas/io/parsers.py in _read(filepath_or_buffer, kwds)
447
448     # Create the parser.
--> 449     parser = TextFileReader(filepath_or_buffer, **kwds)
450
451     if chunksize or iterator:

~/anaconda3/lib/python3.6/site-packages/pandas/io/parsers.py in __init__(self, f, engine, kwds)
816         self.options['has_index_names'] = kwds['has_index_names']
817
--> 818         self._make_engine(self.engine)
819
820     def close(self):

~/anaconda3/lib/python3.6/site-packages/pandas/io/parsers.py in _make_engine(self, engine)
1047     def _make_engine(self, engine='c'):
1048         if engine == 'c':
-> 1049             self._engine = CParserWrapper(self.f, **self.options)
1050         else:
1051             if engine == 'python':

~/anaconda3/lib/python3.6/site-packages/pandas/io/parsers.py in __init__(self, src, **kwds)
1693         kwds['allow_leading_cols'] = self.index_col is not False
1694
-> 1695         self._reader = parsers.TextReader(src, **kwds)
1696
1697         # XXX

pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader.__cinit__()

pandas/_libs/parsers.pyx in pandas._libs.parsers.TextReader._setup_parser_source()

FileNotFoundError: File b'titanic.csv' does not exist

```

Here we see the head of our dataframe. A couple of questions come

to mind. Which variables play a role by determining the probability of surviving the Titanic. Sex and Cabin are not numeric values. How do we convert these to numeric values?

### 3 A first look at the dataset

---

1 data.head()

---

```
<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>PassengerId</th> <th>Survived</th> <th>Pclass</th> <th>Name</th>
<th>Sex</th> <th>Age</th> <th>SibSp</th> <th>Parch</th> <th>Ticket</th>
<th>Fare</th> <th>Cabin</th> <th>Embarked</th> </tr> </thead>
<tbody> <tr> <th>0</th> <td>1</td> <td>0</td> <td>3</td>
<td>Braund, Mr. Owen Harris</td> <td>male</td> <td>22.0</td>
<td>1</td> <td>0</td> <td>A/5 21171</td> <td>7.2500</td> <td>NaN</td>
<td>S</td> </tr> <tr> <th>1</th> <td>2</td> <td>1</td> <td>1</td>
<td>Cumings, Mrs. John Bradley (Florence Briggs Th... </td> <td>female</td>
<td>38.0</td> <td>1</td> <td>0</td> <td>PC 17599</td> <td>71.2833</td>
<td>C85</td> <td>C</td> </tr> <tr> <th>2</th> <td>3</td>
<td>1</td> <td>3</td> <td>Heikkinen, Miss. Laina</td> <td>female</td>
<td>26.0</td> <td>0</td> <td>0</td> <td>STON/O2. 3101282</td>
<td>7.9250</td> <td>NaN</td> <td>S</td> </tr> <tr> <th>3</th>
<td>4</td> <td>1</td> <td>1</td> <td>Futrelle, Mrs. Jacques Heath
(Lily May Peel)</td> <td>female</td> <td>35.0</td> <td>1</td>
<td>0</td> <td>113803</td> <td>53.1000</td> <td>C123</td> <td>S</td>
</tr> <tr> <th>4</th> <td>5</td> <td>0</td> <td>3</td> <td>Allen,
Mr. William Henry</td> <td>male</td> <td>35.0</td> <td>0</td>
<td>0</td> <td>373450</td> <td>8.0500</td> <td>NaN</td> <td>S</td>
</tr> </tbody> </table> </div>
```

Here are the summary statistics of the dataframe. The mean, standard-deviation etc are given in this table.

---

1 data.describe()

---

```

<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
    .dataframe tbody tr th { vertical-align: top; }
    .dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>PassengerId</th> <th>Survived</th> <th>Pclass</th> <th>Age</th>
<th>SibSp</th> <th>Parch</th> <th>Fare</th> </tr> </thead>
<tbody> <tr> <th>count</th> <td>891.000000</td> <td>891.000000</td>
<td>891.000000</td> <td>714.000000</td> <td>891.000000</td> <td>891.000000</td>
<td>891.000000</td> </tr> <tr> <th>mean</th> <td>446.000000</td>
<td>0.383838</td> <td>2.308642</td> <td>29.699118</td> <td>0.523008</td>
<td>0.381594</td> <td>32.204208</td> </tr> <tr> <th>std</th>
<td>257.353842</td> <td>0.486592</td> <td>0.836071</td> <td>14.526497</td>
<td>1.102743</td> <td>0.806057</td> <td>49.693429</td> </tr> <tr>
<th>min</th> <td>1.000000</td> <td>0.000000</td> <td>1.000000</td>
<td>0.420000</td> <td>0.000000</td> <td>0.000000</td> <td>0.000000</td>
</tr> <tr> <th>25%</th> <td>223.500000</td> <td>0.000000</td>
<td>2.000000</td> <td>20.125000</td> <td>0.000000</td> <td>0.000000</td>
<td>7.910400</td> </tr> <tr> <th>50%</th> <td>446.000000</td>
<td>0.000000</td> <td>3.000000</td> <td>28.000000</td> <td>0.000000</td>
<td>0.000000</td> <td>14.454200</td> </tr> <tr> <th>75%</th>
<td>668.500000</td> <td>1.000000</td> <td>3.000000</td> <td>38.000000</td>
<td>1.000000</td> <td>0.000000</td> <td>31.000000</td> </tr> <tr>
<th>max</th> <td>891.000000</td> <td>1.000000</td> <td>3.000000</td>
<td>80.000000</td> <td>8.000000</td> <td>6.000000</td> <td>512.329200</td>
</tr> </tbody> </table> </div>

```

It is possible to search for particular passenger in the dataset. Such as passengers who were older than eighty years.

---

```
1 data[data.Age == 80]
```

---

```

<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
    .dataframe tbody tr th { vertical-align: top; }
    .dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>PassengerId</th> <th>Survived</th> <th>Pclass</th> <th>Name</th>
<th>Sex</th> <th>Age</th> <th>SibSp</th> <th>Parch</th> <th>Ticket</th>
<th>Fare</th> <th>Cabin</th> <th>Embarked</th> </tr> </thead>

```

```
<tbody> <tr> <th>630</th> <td>631</td> <td>1</td> <td>1</td>
<td>Barkworth, Mr. Algernon Henry Wilson</td> <td>male</td> <td>80.0</td>
<td>0</td> <td>0</td> <td>27042</td> <td>30.0</td> <td>A23</td>
<td>S</td> </tr> </tbody> </table> </div>
```

## 4 First figures

To get a good impression of the dataset and the influence of the variables, a couple of diagrams are made using `matplotlib`.

```
1 plt.scatter(data.Age,data.Survived)
2 plt.xlabel('Age')
3 plt.ylabel('Survived')
```

\_\_11\_\_1.png

output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d

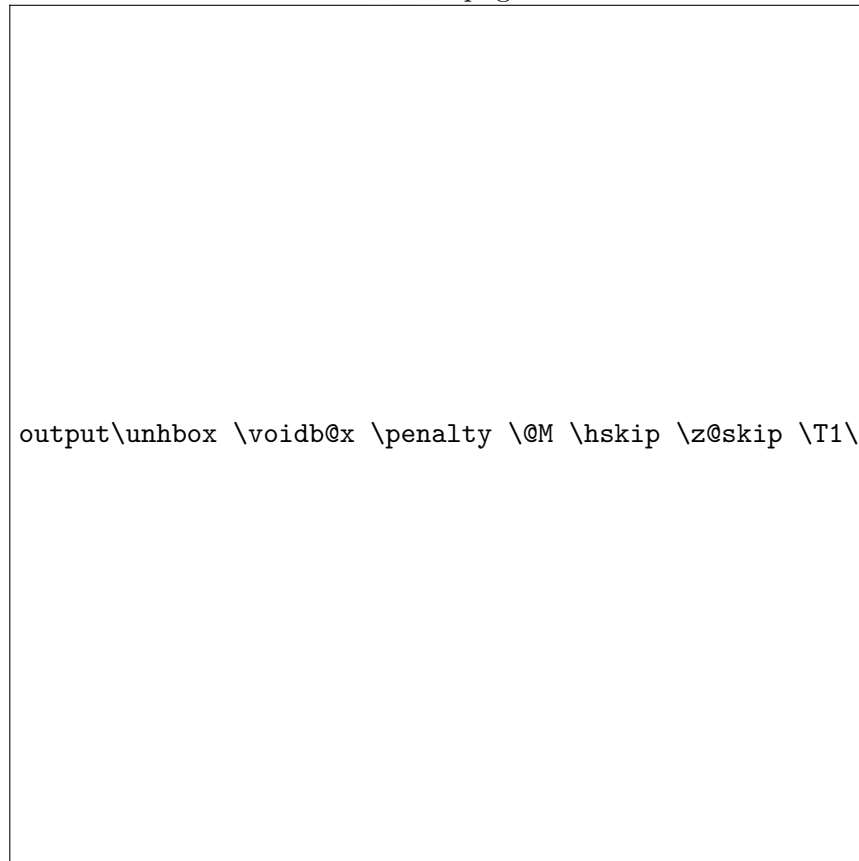
Scatterplots are not always the best choice to illustrate some of the variables. There is not much to say about the variance because of the fact that a lot of points are close to each other. A couple of values however stand out. We see that a passenger or more passengers travelling first class have paid more than 500 pounds for their ticketprice.

---

```
1 plt.scatter(data.Pclass,data.Fare)
2 plt.xlabel('Pclass')
3 plt.ylabel('Fare')
```

---

\_13\_1.png

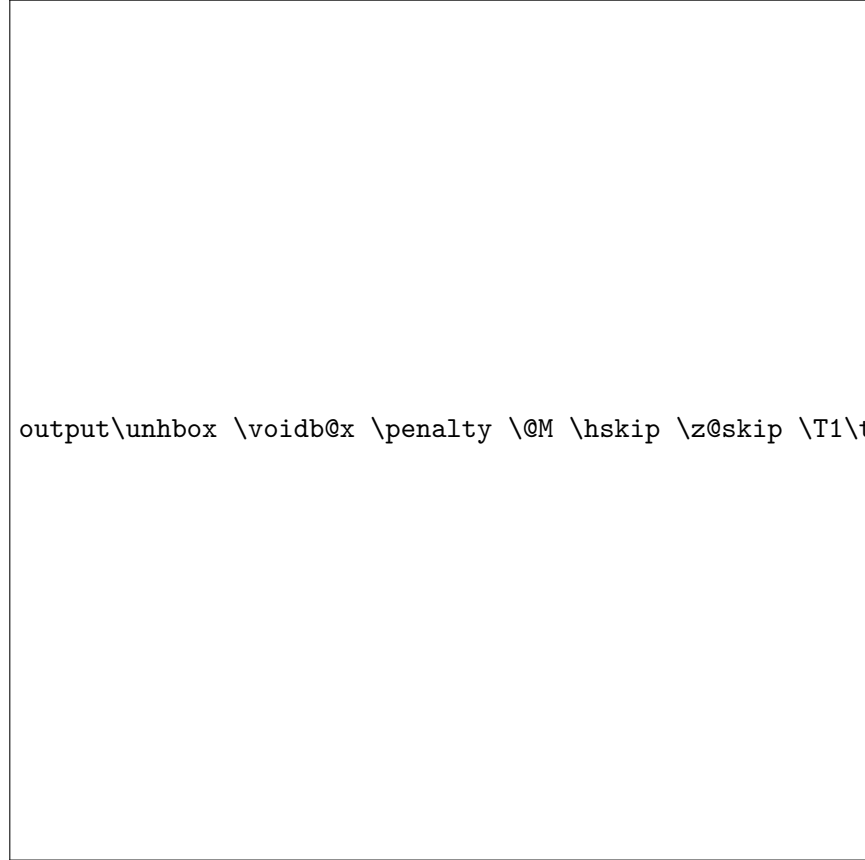


---

```
1 plt.scatter(data.Fare, data.Survived)
2 plt.xlabel('Fare')
3 plt.ylabel('Survived')
```

---

\_14\_1.png



---

```
1 plt.scatter(data.Fare, data.Age)
2 plt.xlabel('Fare')
3 plt.ylabel('Age')
```

---

\_15\_1.png

```
output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d
```

I was curious to see who had paid more than 400 pounds for their ticket. We see that it is easy to make a selection in our dataset using the > sign

```
1 data[data.Fare > 400]
```

```
<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>PassengerId</th> <th>Survived</th> <th>Pclass</th> <th>Name</th>
<th>Sex</th> <th>Age</th> <th>SibSp</th> <th>Parch</th> <th>Ticket</th>
<th>Fare</th> <th>Cabin</th> <th>Embarked</th> </tr> </thead>
<tbody> <tr> <th>258</th> <td>259</td> <td>1</td> <td>1</td>
<td>Ward, Miss. Anna</td> <td>female</td> <td>35.0</td> <td>0</td>
```



```

<td>0</td> <td>PC 17755</td> <td>512.3292</td> <td>NaN</td>
<td>C</td> </tr> <tr> <th>679</th> <td>680</td> <td>1</td>
<td>1</td> <td>Cardeza, Mr. Thomas Drake Martinez</td> <td>male</td>
<td>36.0</td> <td>0</td> <td>1</td> <td>PC 17755</td> <td>512.3292</td>
<td>B51 B53 B55</td> <td>C</td> </tr> <tr> <th>737</th> <td>738</td>
<td>1</td> <td>1</td> <td>Lesurer, Mr. Gustave J</td> <td>male</td>
<td>35.0</td> <td>0</td> <td>0</td> <td>PC 17755</td> <td>512.3292</td>
<td>B101</td> <td>C</td> </tr> </tbody> </table> </div>

```

---

```

1 df_cleaned = data.dropna()
2 df_cleaned['male_dummy'] = (df_cleaned.Sex == 'male') #nieuwe kolom definiëren om male te veranderen in een boolean
3 X = df_cleaned[['Age', 'male_dummy', 'Pclass', 'SibSp', 'Fare']]
4 y = df_cleaned[['Survived']]

```

---

Here we see that we clean our dataset for the first time to make it more suitable for the packages we will be using. All rows with missing values (these are called NaNs, short for Not a Number) are deleted for scikit\\_-learn can't work with NaNs by using `.dropna()`. There are other ways than deleting rows to handle this problem. Replace the NaNs with the mean or to interpolate for example. However the choice was made to delete these rows. Furthermore we see that the problem of the **Sex** column not being a numeric value is handled. The values in the **Sex** column are changed into a boolean. Males are given a **True** and the females are given a **False**. Next a couple of variables have added to **X**. **Age**, **male\_dummy**, **Pclass**, **SibSp**, **Fare** are all numeric values and therefore easy to use.

Here we see the cleaned dataframe with the new added column **male\\_dummy**

---

```

1 df_cleaned.head()

```

---

```

<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>PassengerId</th> <th>Survived</th> <th>Pclass</th> <th>Name</th>
<th>Sex</th> <th>Age</th> <th>SibSp</th> <th>Parch</th> <th>Ticket</th>
<th>Fare</th> <th>Cabin</th> <th>Embarked</th> <th>male_dummy</th>
</tr> </thead> <tbody> <tr> <th>1</th> <td>2</td> <td>1</td>
<td>1</td> <td>Cumings, Mrs. John Bradley (Florence Briggs Th. ...</td>

```

```

<td>female</td> <td>38.0</td> <td>1</td> <td>0</td> <td>PC
17599</td> <td>71.2833</td> <td>C85</td> <td>C</td> <td>False</td>
</tr> <tr> <th>3</th> <td>4</td> <td>1</td> <td>1</td> <td>Futrelle,
Mrs. Jacques Heath (Lily May Peel)</td> <td>female</td> <td>35.0</td>
<td>1</td> <td>0</td> <td>113803</td> <td>53.1000</td> <td>C123</td>
<td>S</td> <td>False</td> </tr> <tr> <th>6</th> <td>7</td>
<td>0</td> <td>1</td> <td>McCarthy, Mr. Timothy J</td> <td>male</td>
<td>54.0</td> <td>0</td> <td>0</td> <td>17463</td> <td>51.8625</td>
<td>E46</td> <td>S</td> <td>True</td> </tr> <tr> <th>10</th>
<td>11</td> <td>1</td> <td>3</td> <td>Sandstrom, Miss. Mar-
guerite Rut</td> <td>female</td> <td>4.0</td> <td>1</td> <td>1</td>
<td>PP 9549</td> <td>16.7000</td> <td>G6</td> <td>S</td> <td>False</td>
</tr> <tr> <th>11</th> <td>12</td> <td>1</td> <td>1</td>
<td>Bonnell, Miss. Elizabeth</td> <td>female</td> <td>58.0</td>
<td>0</td> <td>0</td> <td>113783</td> <td>26.5500</td> <td>C103</td>
<td>S</td> <td>False</td> </tr> </tbody> </table> </div>

```

```

1  from sklearn.linear_model import LogisticRegression
2  logreg = LogisticRegression()
3  logreg.fit(X, y)
4  y_pred = logreg.predict(X)

```

Here we initialize the first regression called logistic regression. We don't split our dataframe in test and training set yet. For a general indication we only use the regressor and fit it on the cleaned dataset. After that we predict on the same dataset.

```

1  logreg.coef_

```

Here we see the outcome of our first try with the logistic regression. To interpret these coefficients, let's look at the order of the columns in X:

```

1  X.head()

```

```

<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-
align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>

```

```

<th>Age</th> <th>male_dummy</th> <th>Pclass</th> <th>SibSp</th>
<th>Fare</th> </tr> </thead> <tbody> <tr> <th>1</th> <td>38.0</td>
<td>False</td> <td>1</td> <td>1</td> <td>71.2833</td> </tr>
<tr> <th>3</th> <td>35.0</td> <td>False</td> <td>1</td> <td>1</td>
<td>53.1000</td> </tr> <tr> <th>6</th> <td>54.0</td> <td>True</td>
<td>1</td> <td>0</td> <td>51.8625</td> </tr> <tr> <th>10</th>
<td>4.0</td> <td>False</td> <td>3</td> <td>1</td> <td>16.7000</td>
</tr> <tr> <th>11</th> <td>58.0</td> <td>False</td> <td>1</td>
<td>0</td> <td>26.5500</td> </tr> </tbody> </table> </div>

```

## 5 Graphic illustration of a prediction

One of the first graphic illustrations of the relation between fare, age and survival. The relation is not very clear but we see that the higher the fare the more people survived and the higher the age the less people survived. However, this figure is not very accurate, because of the fact that only three variables were used.

---

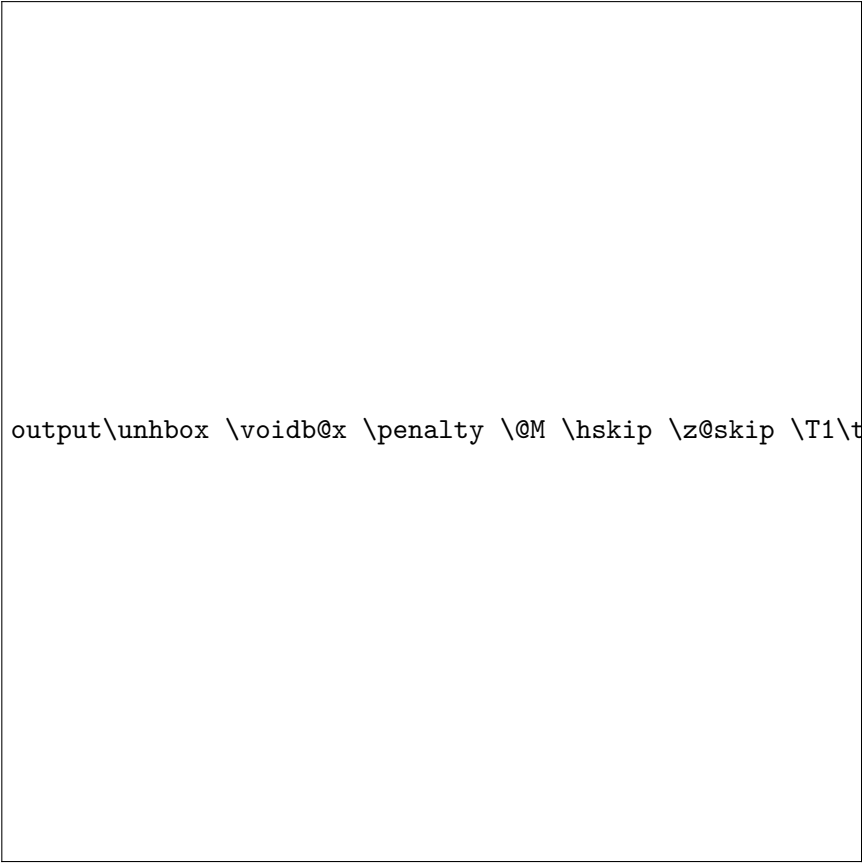
```

1     survived = df_cleaned[df_cleaned.Survived == 1]
2     not_survived = df_cleaned[df_cleaned.Survived == 0]
3
4     plt.scatter(survived.Fare, survived.Age, marker='^', label = 'survived')
5     plt.scatter(not_survived.Fare, not_survived.Age, marker='^', label = 'not survived')
6     plt.legend()

```

---

\_\_28\_\_1.png



output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d

---

```
1 P = df_cleaned[['Pclass', 'Fare', 'Age', 'male_dummy']]
```

---

We select from our `df_cleaned` only the columns with numeric values. This is convenient for the splitting into train and testsets, for `scikit-learn` can only work with numbers. Difference between `P` and `X` here is that `X` also has the column `siblings`, whereas `P` only has four columns

---

```
1 P.head()
```

---

```
<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
```

Pclass	Fare	Age	male_dummy
1	71.2833	38.0	False
3	53.1000	35.0	False
6	51.8625	54.0	True
10	3	16.7000	4.0
11	1	26.5500	58.0

---

```

1 from sklearn.linear_model import LogisticRegression
2 logreg = LogisticRegression(fit_intercept=True)
3 logreg.fit(P, y)
4 y_pred = logreg.predict(P)

```

---

We fit our regressor on our dataset and predict on that same dataset. Once again without splitting into train and testset. Just to get a general idea about the values of the coefficients.

---

```

1 logreg.coef_

```

---

One could interpret the found coefficients as follows: The coefficients for class and fare are positive, which may indicate that the higher the class and price paid for a ticket, the higher the chance of surviving the Titanic. When we look at age and sex we see the exact opposite for the coefficients are negative. The higher the age the lower your chances and if you were a man on board of the titanic your chances of surviving were lower.

---

```

1 from sklearn.neighbors import KNeighborsClassifier

```

---

Another regression is used in the following lines. (explanation K nearest neighbours)

---

```

1 knn = KNeighborsClassifier(n_neighbors=6)

```

---



---

```

1 knn.fit(P,y)

```

---



---

```

1 prediction = knn.predict(P)

```

---

---

```
1 P.shape
```

---

---

```
1 print('Prediction{}'.format(prediction))
```

---

Here we see one of our first predictions. 1 indicates the passenger has survived and 0 indicates that the passenger has died

Elke persoon heeft andere karakteristieken, dus dit zijn voorspellingen per persoon. Dus er komt een kans uit en dan kijkt de regressor, boven of onder 0.5

---

```
1 knn.score(P,y)
```

---

This score gives a number between 0 and 1 and gives an impression of the accuracy of our model. However, this accuracy is not an indication of how well our model performs (explanation spam mail etc.)

---

```
1 P.head()
```

---

```
<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>Pclass</th> <th>Fare</th> <th>Age</th> <th>male_dummy</th>
</tr> </thead> <tbody> <tr> <th>1</th> <td>1</td> <td>71.2833</td>
<td>38.0</td> <td>False</td> </tr> <tr> <th>3</th> <td>1</td>
<td>53.1000</td> <td>35.0</td> <td>False</td> </tr> <tr> <th>6</th>
<td>1</td> <td>51.8625</td> <td>54.0</td> <td>True</td> </tr>
<tr> <th>10</th> <td>3</td> <td>16.7000</td> <td>4.0</td> <td>False</td>
</tr> <tr> <th>11</th> <td>1</td> <td>26.5500</td> <td>58.0</td>
<td>False</td> </tr> </tbody> </table> </div>
```

---

```
1 q = df_cleaned.Survived
```

---

---

```
1 q.head()
```

---

---

```
1 from sklearn.linear_model import LogisticRegression
```

---

---

```

1  from sklearn.model_selection import train_test_split

```

---

```

1  logreg = LogisticRegression()
2  P_train, P_test, q_train, q_test = train_test_split(P,q, test_size=0.2, random_state=42)
3  logreg.fit(P_train, q_train)
4  q_pred = logreg.predict(P_test)

```

---

Here we see the dataset being split into a test and a training set. The arguments give us information about how much of our data we use as a test\\_set and how much of our data we use as a training\\_set. This and the parameters will be varied to see which parameter gives the best prediction. We fit our regressor on the training\\_set and predict on the test\\_set.

---

```

1  print('Prediction {}'.format(q_pred))

```

---



---

```

1  P_train.head()

```

---

```

<div> <style scoped> .dataframe tbody tr th:only-of-type { vertical-align: middle; }
.dataframe tbody tr th { vertical-align: top; }
.dataframe thead th { text-align: right; } </style> <table border="1"
class="dataframe"> <thead> <tr style="text-align: right;"> <th></th>
<th>Pclass</th> <th>Fare</th> <th>Age</th> <th>male_dummy</th>
</tr> </thead> <tbody> <tr> <th>331</th> <td>1</td> <td>28.5000</td>
<td>45.5</td> <td>True</td> </tr> <tr> <th>336</th> <td>1</td>
<td>66.6000</td> <td>29.0</td> <td>True</td> </tr> <tr> <th>193</th>
<td>2</td> <td>26.0000</td> <td>3.0</td> <td>True</td> </tr>
<tr> <th>75</th> <td>3</td> <td>7.6500</td> <td>25.0</td> <td>True</td>
</tr> <tr> <th>248</th> <td>1</td> <td>52.5542</td> <td>37.0</td>
<td>True</td> </tr> </tbody> </table> </div>

```

---

```

1  from sklearn.metrics import roc_auc_score
2  q_pred_prob = logreg.predict_proba(P_test)[: ,1]
3  roc_auc_score(q_test, q_pred_prob)

```

---

When the test\\_size is changed from 0.4 to 0.2 , the score increases with more than 10%. This makes sense because a smaller test\\_set gives a higher accuracy score.

---

```

1  from sklearn.model_selection import cross_val_score
2  cv_scores = cross_val_score(logreg, P, q, cv=5, scoring='roc_auc')
3  print(cv_scores)

```

---



---

```

1  len(P)

```

---



---

```

1  len(prediction)

```

---

Seaborn countplot option? plt.figure()

sns.countplot(x='education', hue='party', data=df, palette='RdBu')

plt.xticks([0,1], ['No', 'Yes']) plt.show()

Given all the variables (age, gender, place of boarding etc.), you make a linear function ( $a_1x_1 + a_2x_2 + \dots + a_nx_n + b$ ). Computer puts this in the logistic function for  $x$ . For a particular  $x$ , you get a value between zero and one. This is your chance of survival. Boundary is 0,5.  $X < 0,5$  passenger didn't survive. Computer tries to plot a logistic function where  $R^2$  is as small as possible. This is called the fitting process. The logistic function has to be as close to the datapoints as possible.

---

```

1  sns.countplot?

```

---



---

```

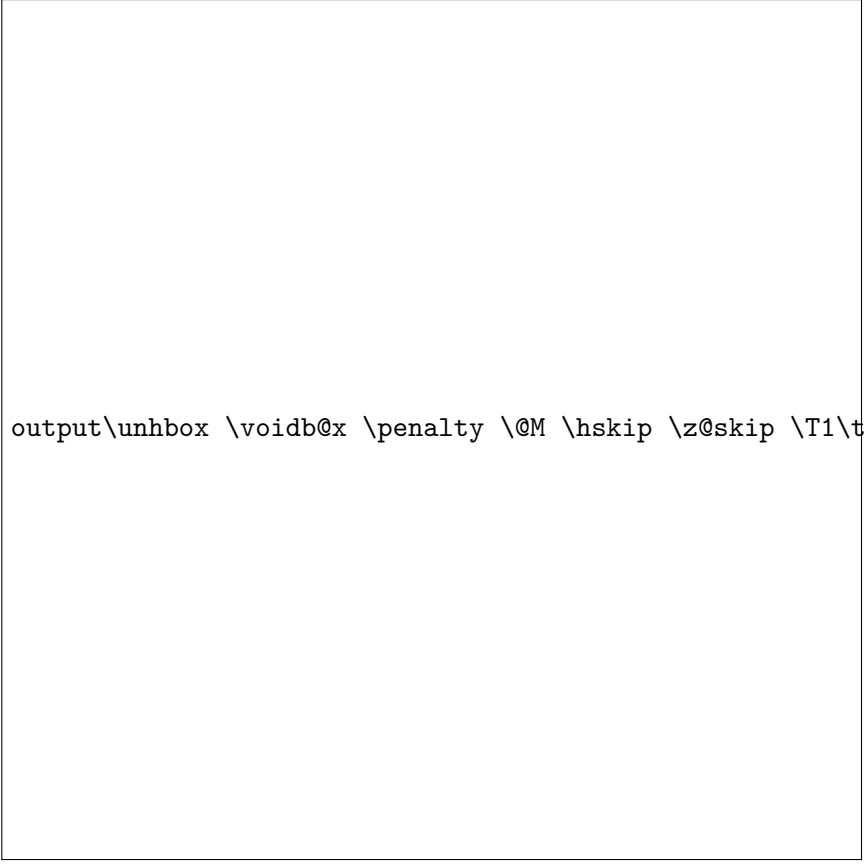
1  sns.set(style="darkgrid")
2  ax = sns.countplot(x="Pclass", hue="Survived", data=data, palette="Set3")

```

---

\_64\_0.png





```
output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d
```

More people in class 3 than in class 1, makes it difficult to compare and draw a conclusion. Percentage? In general, we cannot draw a conclusion regarding survival probabilities because there were more people in class 3 than in one 1. In the third class, more passengers died than survived. In the first class, more people survived than perished. We cannot compare the results from the first class to the third class. The plot only shows us one variable. This is another reason why we cannot be sure about the influence of class on the chance of survival. *Simpson paradox*

---

```
1 sns.set(style="darkgrid")
2 ax = sns.countplot(x="Age",hue="Survived", data=data, palette="Set1")
```

---

\_\_66\_\_0.png

output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d

---

```
1 sns.set(style="darkgrid")
2 ax = sns.countplot(x="Fare",hue="Survived", data=data)
```

---

\_67\_0.png

output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d

---

```
1 sns.set(style="darkgrid")
2 ax = sns.countplot(x="male_dummy",hue="Survived", data=df_cleaned, palette="Set2")
```

---

\_68\_0.png



output\unhbox \voidb@x \penalty \@M \hskip \z@skip \T1\textunderscore \protect \d