There are three files in the data: (1) train.csv, (2) test.csv, and (3) gender_submission.csv.

(1) train.csv train.csv contains the details of a subset of the passengers on board (891 passengers, to be exact -- where each passenger gets a different row in the table). The values in the second column ("Survived") can be used to determine whether each passenger survived or not:

if it's a "1", the passenger survived. if it's a "0", the passenger died.
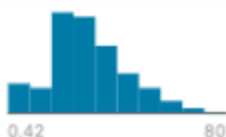
(2) test.csv

< **train.csv** (59.76 KB)

Detail    Compact    Column

**About this file**

contains data

| ⇔ PassengerId | | # Survived | | # Pclass | | A Name | | A Sex | | # Age | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | male 65%<br>female 35% | | | |
| | | | | | | 891<br>unique values | | | | | |
| 1 | 891 | 0 | 1 | 1 | 3 | | | | | 0.42 | 80 |
| 1 | | 0 | | 3 | | Braund, Mr. Owen Harris | | male | | 22 | |
| 2 | | 1 | | 1 | | Cumings, Mrs. John Bradley (Florence Briggs Thayer) | | female | | 38 | |
| 3 | | 1 | | 3 | | Heikkinen, Miss. Laina | | female | | 26 | |

Using the patterns you find in train.csv, you have to predict whether the other 418 passengers on board (in test.csv) survived.

(3) gender_submission.csv The gender_submission.csv file is provided as an example that shows how you should structure your predictions. It predicts that all female passengers survived, and all male passengers died. Your hypotheses regarding survival will probably be different, which will lead to a different submission file. But, just like this file, your submission should have:

a "PassengerId" column containing the IDs of each passenger from test.csv. a "Survived" column (that you will create!) with a "1" for the rows where you think the passenger survived, and a "0" where you predict that the passenger died.

```
# This Python 3 environment comes with many helpful analytics libraries installed
# It is defined by the kaggle/python docker image: https://github.com/kaggle/docker-python
# For example, here's several helpful packages to load in

import numpy as np # linear algebra
import pandas as pd # data processing, CSV file I/O (e.g. pd.read_csv)

# Input data files are available in the "../input/" directory.
# For example, running this (by clicking run or pressing Shift+Enter) will list all files under the

import os
```

```python
for dirname, _, filenames in os.walk('/kaggle/input'):
    for filename in filenames:
        print(os.path.join(dirname, filename))

# Any results you write to the current directory are saved as output.
```

```python
train_data = pd.read_csv("/content/train.csv")
train_data.head()
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 |

```python
test_data = pd.read_csv("/content/train.csv")
test_data.head()
```

| | PassengerId | Survived | Pclass | Name | Sex | Age | SibSp | Parch | Ticket | Fare | Cabin |
|---|---|---|---|---|---|---|---|---|---|---|---|
| **0** | 1 | 0 | 3 | Braund, Mr. Owen Harris | male | 22.0 | 1 | 0 | A/5 21171 | 7.2500 | NaN |
| **1** | 2 | 1 | 1 | Cumings, Mrs. John Bradley (Florence | female | 38.0 | 1 | 0 | PC 17599 | 71.2833 | C85 |

```python
women = train_data.loc[train_data.Sex == 'female']["Survived"]
rate_women = sum(women)/len(women)

print("% of women who survived:", rate_women)
```

```
% of women who survived: 0.7420382165605095
```

```python
men = train_data.loc[train_data.Sex == 'male']["Survived"]
rate_men = sum(men)/len(men)

print("% of men who survived:", rate_men)
```

```
% of men who survived: 0.18890814558058924
```

```python
from sklearn.ensemble import RandomForestClassifier

y = train_data["Survived"]

features = ["Pclass", "Sex", "SibSp", "Parch"]
X = pd.get_dummies(train_data[features])
```

```
X_test = pd.get_dummies(test_data[features])

model = RandomForestClassifier(n_estimators=100, max_depth=5, random_state=1)
model.fit(X, y)
predictions = model.predict(X_test)

output = pd.DataFrame({'PassengerId': test_data.PassengerId, 'Survived': predictions})
output.to_csv('submission.csv', index=False)
print("Your submission was successfully saved!")
```

```
Your submission was successfully saved!
```