

Overview

The data has been split into two groups:

```
training set (train.csv)
test set (test.csv)
```

The training set should be used to build your machine learning models. For the training set, we provide the outcome (also known as the “ground truth”) for each passenger. Your model will be based on “features” like passengers’ gender and class. You can also use feature engineering to create new features.

The test set should be used to see how well your model performs on unseen data. For the test set, we do not provide the ground truth for each passenger. It is your job to predict these outcomes. For each passenger in the test set, use the model you trained to predict whether or not they survived the sinking of the Titanic.

Data Dictionary

Variable	Definition	Key
survival	Survival	0 = No, 1 = Yes
pclass	Ticket class	1 = 1st, 2 = 2nd, 3 = 3rd
sex	Sex	

Age Age in years
sibsp # of siblings / spouses aboard
the Titanic
parch # of parents / children aboard
the Titanic
ticket Ticket number
fare Passenger fare
cabin Cabin number
embarked Port of Embarkation C
= Cherbourg, Q = Queenstown, S =
Southampton
Variable Notes

pclass: A proxy for socio-economic status
(SES)
1st = Upper
2nd = Middle
3rd = Lower

age: Age is fractional if less than 1. If
the age is estimated, is it in the form
of xx.5

sibsp: The dataset defines family
relations in this way...
Sibling = brother, sister, stepbrother,
stepsister
Spouse = husband, wife (mistresses and
fiancés were ignored)

parch: The dataset defines family
relations in this way...
Parent = mother, father

Child = daughter, son, stepdaughter,
stepson

Some children travelled only with a
nanny, therefore parch=0 for them.