

Titanic - Machine Learning from Disaster

Start here! Predict survival on the Titanic and get familiar with ML basics



Titanic - Machine Learning from Disaster

[**Submit Prediction**](#)

Dataset Description

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.

We also include **gender_submission.csv**, a set of predictions that assume all and only female passengers survive, as an example of what a submission file should look like.

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	
embark ed	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.

Files

3 files

Size

93.08 kB

Type

CSV

License

Subject to Competition Rules

gender_submission.csv(3.26 kB)

About this file

An example of what a submission file should look like.

These predictions assume only female passengers survive.

PassengerId

Survived

Label	Coun t
-------	-----------

892.00

- 42

933.70

933.70

- 42

975.40

975.40

- 42

1017.1

0 - 41

1058.8

0

1058.8

0 - 42

1100.50

1100.50
- 42
1142.20

1142.20
- 41
1183.90

1183.90
-
1225.6 42
0

1225.6
0 -
1267.3 42
0

1267.3
0 -
1309.0 42
0

892

1309

Label Count
I t

0.00
- 266
0.10

0.90
- 152
1.00

0

1

89208931894089508961897089818990900190109020903090419050906190719080909091
01911191209130914191509161917091819190920092109220923092419251926092709281
9291930093109320933093409351936193709380939094019411

Data Explorer

93.08 kB

gender_submission.csv

test.csv

train.csv

Summary

3 files

25 columns

[Download All](#)

[Download data](#)

Metadata

License

Subject to Competition Rules