

Pune Institute of Computer Technology Dhankawadi, Pune

A MINI PROJECT REPORT ON

Who Survived the Titanic shipwreck prediction using Machine Learning

**SUBMITTED TO THE SAVITRIBAI PHULE PUNE UNIVERSITY, PUNE IN THE PARTIAL
FULFILLMENT OF THE LP III
BE (COMPUTER ENGINEERING)**

SUBMITTED BY

Roll No.	Name
41471	Prachi Said

Under the Guidance of

Prof. S. P. Shintre



**DEPARTMENT OF COMPUTER ENGINEERING
Academic Year 2024-25**



**PUNE INSTITUTE OF COMPUTER TECHNOLOGY DEPARTMENT OF COMPUTER
ENGINEERING**

CERTIFICATE

This is to certify that the Mini Project report of LP III (ML) entitled
"Who survived the titanic shipwreck prediction using machine learning"

Submitted by

Roll No.	Name
41471	Prachi Said

has satisfactorily completed a micro project report under the guidance of Prof. S.
P. Shintre towards the partial fulfillment of BE Computer Engineering, Academic Year
2024-25 of Savitribai Phule Pune University.

Prof. S. P. Shintre
(Lab Guide)

Dr. Geetanjali Kale
(H. O. D)

Date: / / Place:
Pune

Introduction

The goal of the project was to predict the survival of passengers based off a set of data. We used Kaggle competition "Titanic: Machine Learning from Disaster" (see <https://www.kaggle.com/c/titanic/data>) to retrieve necessary data and evaluate accuracy of our predictions. The historical data has been split into two groups, a 'training set' and a 'test set'. For the training set, we are provided with the outcome (whether or not a passenger survived). We used this set to build our model to generate predictions for the test set.

For each passenger in the test set, we had to predict whether or not they survived the sinking. Our score was the percentage of correctly predictions.

Training and Test Data

Training and Test data come in CSV file and contain the following fields:

- Passenger ID
- Passenger Class
- Name
- Sex
- Age
- Number of passenger's siblings and spouses on board
- Number of passenger's parents and children on board
- Ticket
- Fare
- Cabin
- City where passenger embarked

PassengerId	Survived	Pclass	Name	Sex	Age	SibSp	Parch	Ticket	Fare	Cabin	Embarked
1	0	3	"Braund, Mr. Owen Harris"	male	22	1	0	A/5 21171	7.25		S
2	1	1	"Cumings, Mrs. John Bradley (Florence Briggs Thayer)"	female	38	1	0	PC 17599	71.2833	C85	C
3	1	3	"Heikkinen, Miss. Laina"	female	26	0	0	STON/O2. 3101282	7.925		S
4	1	1	"Futrelle, Mrs. Jacques Heath (Lily May Peel)"	female	35	1	0	113803	53.1	C123	S
5	0	3	"Allen, Mr. William Henry"	male	35	0	0	373450	8.05		S
6	0	3	"Moran, Mr. James"	male		0	0	330877	8.4583		Q
7	0	1	"McCarthy, Mr. Timothy J"	male	54	0	0	17463	51.8625	E46	S
8	0	3	"Palsson, Master. Gosta Leonard"	male	2	3	1	349909	21.075		S
9	1	3	"Johnson, Mrs. Oscar W (Elisabeth Vilhelmina Berg)"	female	27	0	2	347742	11.1333		S
10	1	2	"Nasser, Mrs. Nicholas (Adele Achem)"	female	14	1	0	237736	30.0708		C
11	1	3	"Sandstrom, Miss. Marguerite Rut"	female	4	1	1	PP 9549	16.7	G6	S
12	1	1	"Bonnell, Miss. Elizabeth"	female	58	0	0	113783	26.55	C103	S
13	0	3	"Saunderscock, Mr. William Henry"	male	20	0	0	A/5. 2151	8.05		S
14	0	3	"Andersson, Mr. Anders Johan"	male	39	1	5	347082	31.275		S
15	0	3	"Vestrom, Miss. Hulda Amanda Adolfina"	female	14	0	0	350406	7.8542		S
16	1	2	"Hewlett, Mrs. (Mary D Kingcome) "	female	55	0	0	248706	16		S
17	0	3	"Rice, Master. Eugene"	male	2	4	1	382652	29.125		Q
18	1	2	"Williams, Mr. Charles Eugene"	male		0	0	244373	13		S
19	0	3	"Vander Planke, Mrs. Julius (Emelia Maria Vandemoortele)"	female	31	1	0	345763	18		S
20	1	3	"Masselmani, Mrs. Fatima"	female		0	0	2649	7.225		C

Feature Engineering

Since the data can have missing fields, incomplete fields, or fields containing hidden information, a crucial step in building any prediction system is Feature Engineering. For instance, the fields Age, Fare, and Embarked in the training and test data, had missing values that had to be filled in. The field Name while being useless itself, contained passenger's Title (Mr., Mrs., etc.), we also used passenger's surname to distinguish families on board of Titanic. Below is the list of all changes that has been made to the data.

Extracting Title from Name

The field Name in the training and test data has the form "Braund, Mr. Owen Harris". Since name is unique for each passenger, it is not useful for our prediction system. However, a passenger's title can be extracted from his or her name. We found 10 titles:

Index	Title	Number of occurrences
0	Col.	4
1	Dr.	8
2	Lady	4
3	Master	61
4	Miss	262
5	Mr.	757
6	Mrs.	198
7	Ms.	2
8	Rev.	8
9	Sir	5

Calculating Family Size

It seems advantageous to calculate family size as follows

$\text{Family_Size} = \text{Parents_Children} + \text{Siblings_Spouses} + 1$.

Extracting Deck from Cabin

The field Cabin in the training and test data has the form "C85", "C125", where C refers to the deck label. We found 8 deck labels: A, B, C, D, E, F, G, T. We see deck label as a refinement of the passenger's class field since the decks A and B were intended for passengers of the first class, etc.

Extracting Ticket_Code from Ticket

The field Ticket in the training and test data has the form "A/5 21171". Although we couldn't

understand meaning of letters in front of numbers in the field Ticket, we extracted those letters and used them in our prediction system. We found the following letters :

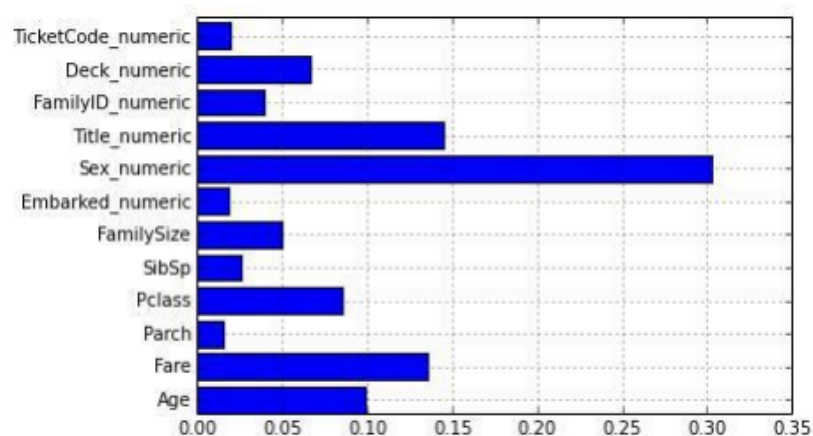
Index	Ticket Code	Number of occurrences
0	No Code	961
1	A	42
2	C	77
3	F	13
4	L	1
5	P	98
6	S	98
7	W	19

Filling in missing values in the fields Fare, Embarked, and Age

Since the number of missing values was small, we used median of all Fare values to fill in missing Fare fields, and the letter 'S' (most frequent value) for the field Embarked. In the training and test data, there was significant amount of missing Ages. To fill in those, we used Linear Regression algorithm to predict Ages based on all other fields except Passenger_ID and Survived.

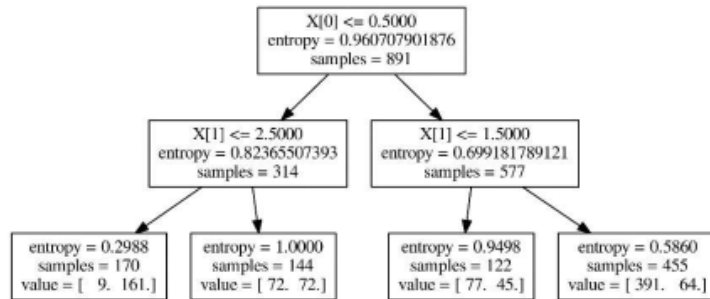
Importance of fields

Decision Trees algorithm in the library SciKit-Learn allows to evaluate importance of each field used for prediction. Below is the chart displaying importance of each field.



Decision Trees

Our prediction system is based on growing Decision Trees to predict the survival status. A typical Decision Tree is pictured below

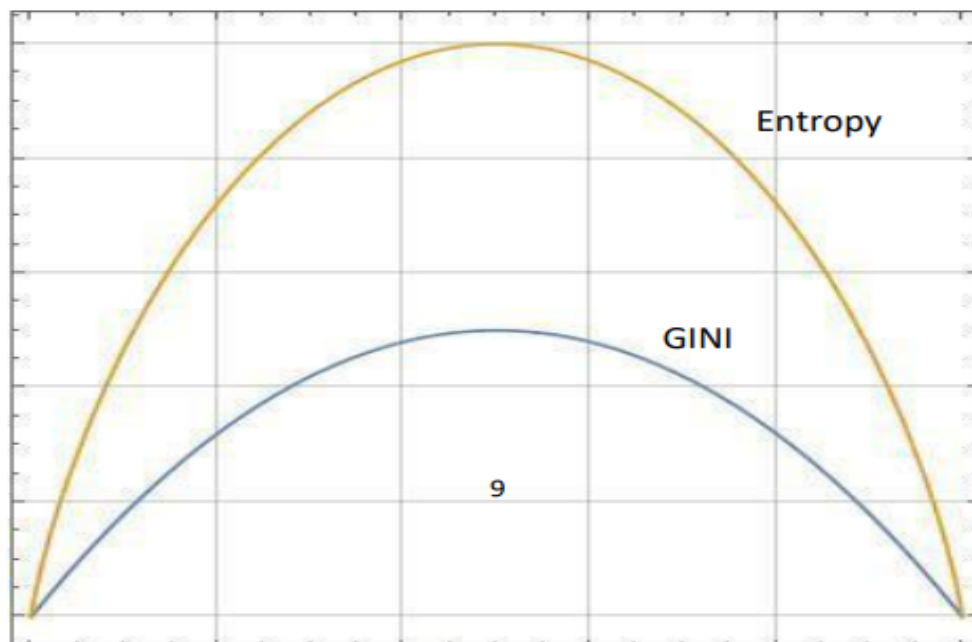


The basic algorithm for growing Decision Tree:

1. Start at the root node as parent node.
2. Split the parent node based on field $X[i]$ to minimize the sum of child nodes uncertainty (maximize information gain).
3. Assign training samples to new child nodes.
4. Stop if leave nodes are pure or early stopping criteria is satisfied, otherwise repeat step 1 and 2 for each new child node.

Stopping Rules:

1. The leaf nodes are pure
2. A maximal node depth is reached
3. Splitting a node does not lead to an information gain



Conclusion

As a result of our work, we gained valuable experience of building prediction systems and achieved our best score on Kaggle: 80.383% of correct predictions (in Kaggle leaderboard, it corresponds to positions 477 - 881 out of 3911 participants).

We performed featured engineering techniques :

- Changed alphabetic values to numeric
- Calculated family size
- Extracted title from name and deck label from ticket number
- Used linear regression algorithm to fill in missing ages

We used several prediction algorithms in python

- Decision tree
- Random forests
- Extra trees

We achieved our best score 80.383% correct predictions.