

Introduction

The titanic was catastrophic as two-thirds of her passengers died in her maiden voyage. The list of those that passed has been published and was turned into a data science competition to see who can create a model that could predict the survival rate with the highest accuracy. The question then became who died and who lived and what correlations existed between them. Most people have watched the movie *Titanic* and saw that Rose lived and Jack died. Despite your opinions on the matter, however, using mathematical modeling it will show that even if there was room on the raft, Jack still most likely would've died and Rose would've lived. Now that a probable outcome is proven, how would you fare on the titanic?

Data Sets & Variables

The predictions were created by analyzing two data sets: the training set (Figure 1) which contained 11 variables including the survival, and the testing set (Figure 2) which contained all variables except the survival. The model was created using the training set then evaluated with the test set. For missing data, the variable was replaced or disregarded. RStudio was used for analysis.

Passengerld	Surviva	l Class	Class Name		Sex		Sibling/Spouse	Parent/Child	Ticket	t	Fare Paid	Cabin	Embarked
1	Passed	3	Braund, Mr. Owen Harris		male	22	1	0	A/5 211	71	7.25		S
2	Survived	1 1	Cumings, Mrs. Florence		female	38	1	0	PC 1759	99	71.2833	C85	С
3	Survived	d 3	Heikkinen, Miss. Laina		female	26	0	0	STON/O2. 3101282		7.925		S
4	Survived	1 1	Futrelle, Mrs. Jacques He	ath	female	35	1	0	113803		53.1	C123	S
Figure 1. Training set with known outcome													
Passengerid	Class	Name Sex Age Sibling/Spouse Parent/Child Ticket Fare Cabin Embarked S							Survival				
892	3	Kelly, Mr. James		m	ale 3	34.5	0	0	330911	7.829	2	Q	?
893	3	Wilkes, Mrs. James (Ellen Needs)		fen	nale	47	1	0	363272	7		S	?
894	2	Myles, Mr. Thomas Francis			ale	62	0	0	240276	9.687	5	Q	?
895	3	Wirz, Mr. Albert			ale	27	0	0	315154	8.662	5	S	?

Feature Engineering

More variables were created by using feature engineering — the modification of preexisting variables whether it is pulling apart, cutting apart and/or merging of data. The following variables were created from stories that were prevalent.

A famous maritime sayings is "women and children first". An child variable was created by making every passenger under the age of 18 a child and those above an adult. From Figure 3 it is seen that this holds true for the titanic.

Child	Sex
NO	F
Yes	F
NO	М
Yes	М
T1+	2 T 11 (

Figure 3. Table of survival probabilities for children of both sexes.

Fare2 <10	Survived 19.94048%
10-20	42.45810%
20-30	42.64706%
30+	58.75000%

Figure 4. Table of survival probabilities for the 4 fare buckets

Potentially the most important saying analyzed was that "prestige lives". By pulling apart and cutting up the name, the title of the person was isolated. It was again found that the saying held true that those with a higher title had better survival odds.

survival odds.











Titanic: Surviving the Disaster

Department of Mathematics, University of Wisconsin-La Crosse, WI 54601

Nicole Nelson, Song Chen

Survived
75.28958%
69.09091%
16.57033%
39.65517%

The next saying explored was "it pays to be rich". Fares were merged into four price ranges and it was found that the fare was directly correlated to



Simple Models

To start, the first prediction was that everyone dies. This came back with a 63% accuracy which at the time was better than half the competition.



Next it was observed that about two-thirds of the passengers were male so the prediction only men died was made. This came back with an improved accuracy of 77%.

Since men had a terrible survival rate, women were looked at. It was guessed that third class women would have low survival chances, so it was predicted that men and lower class women would die, but this came back with the exact same accuracy of 77%.

Decision Trees

Instead of examining each variable by hand, decision trees were used to automate the process. Decision trees are a classification model that scans through all the variables and finds the one that causes the largest split (most importance). This method resulted with an 79% accuracy using feature engineered variables. Try reading the tree yourself using Mr. Kelly James. Does he survive?



Random Forest

Due to decision trees splitting on the variable of highest importance, the model is unable to distinguish if the variable should be the first split or later in the tree to reduce error. Random forest creates multiple (1000+ at times) decision trees and akes the average to determine which variable is most important at each split. Random forest came back with an accuracy of 80% (Figure 6).



Figure 7. An example of how decision trees are averaged into a random forest. The first two decision trees had a high accuracy while the third had a low accuracy. Once combined into random forest it was found that title was significant whereas embarked was not



Classification Problem

From analysis it is clear that the question posed is a classification problem. Various methods were analyzed and it was found that conditional inference random forest and logistic regression (Figure 8) proved to be the most effective with an accuracy of 81% and 87% respectively.



Figure 8. ANOVA table which was part of the logistic regression model. The p values close to zero denote that the variable is significant.

Accuracy Measurement

To measure accuracy, models were submitted to Kaggle which was the server hosting the data competition. An immediate score was provided.



Conclusions

Of the classification methods logistic regression proved to be the most effective model in predicting the survival with an accuracy of 87%. Feature engineering was also important as more information was able to be extracted; it was seen that title was the key component in predictions. Following that, age, class, and sex were important factors in the decision.

After the last submission the student was ranked 102 of 6,338 competitors.

> 102 • 14 Nickle.Nelz **Figure 10**. Screen clippings from the competition of the ranking

References

Nov. 2016.

Mick001. "Logistic Regression Tutorial Code." Gist. N.p., 13 Sept. 2015. Web. 03 Nov. 2016. Stephens, Trevor. "Titanic: Getting Started With R." Trevor Stephens. N.p., 09 Jan. 2014. Web. 03 Nov. 2016.

Web. 03 Nov. 2016.

Trevorstephens. "Trevorstephens/titanic." GitHub. N.p., 19 Jan. 2014. Web. 03 Nov. 2016.





			<u> </u>			
	Df	Deviance	Resid. Df	Resid.	Dev Pr(>Chi)	
			799	1066.3	3	
	1	84.779	798	981.5	6 < 2.2e-16	***
	1	240.084	797	741.4	7 < 2.2e-16	***
	1	20.638	796	720.8	3 5.549e-06	***
	1	13.549	795	707.2	8 0.0002325	***
	1	0.854	794	706.4	3 0.3554687	
	1	0.736	793	705.70	0.3910806	
	2	1.435	791	704.2	6 0.4879177	
	3	2.974	788	701.2	9 0.3956678	
	10	45.596	5 77	8 655.6	9 1.698e-06	***
ze	0	0.000	778	655.6	9	
2	21	54.478	3 75	7 01.2	1 8.419e-05	***

ng Started Prediction Competition									
ic: Machine Learning from Disaster									
re! Predict survival on the Titanic and get familiar with ML basics									
Kaggle · 6,385 teams · 3 years to go									
Overview	Data	Kernels	Discussion	Leaderboard	More	My Submissions	Submit Predictions		
Complete									
Your submission scored 0.81340.									
Figure 9	. Scre	en clipp	ping of th	e competit	ion as v	well as the subm	nission		

Alice, Michy. "How to Perform a Logistic Regression in R." Rbloggers. N.p., 13 Sept. 2015. Web. 03

"Titanic: Machine Learning from Disaster." Kaggle: Your Home for Data Science. N.p., 28 Sept. 2012.

Dr. Song Chen & Dr. Chad Vidden Colin Giles & Elizabeth McMahon For more information contact Dr. Chen at schen@uwlax.edu or Dr. Vidden at cvidden@uwlax.edu