

Modeling Titanic Survival

Qiushi Yan

^aBeijing, China

ARTICLE HISTORY

Compiled October 25, 2020

ABSTRACT

This case study showcases the development of a binary logistic model to predict the possibility of survival in the loss of Titanic. I demonstrate the overall modeling process, including preprocessing, exploratory analysis, model fitting, adjustment, bootstrap validation and interpretation as well as other relevant techniques such as redundancy analysis and multiple imputation for missing data. The motivation and justification behind critical statistical decisions are explained. This analysis is fully reproducible with all source R code and text.

<http://www.crema-research.ch/papers/2009-03.pdf>

Who Survived Titanic? A Logistic Regression Analysis: <https://sci-hub.do/>
<https://journals.sagepub.com/doi/pdf/10.1177/084387140401600205>

[https://www.insider.com/titanic-secrets-facts-2018-4#](https://www.insider.com/titanic-secrets-facts-2018-4#at-the-memorial-of-frederick-fleet-one-of-the-lookouts-a-prankster-left-a-pair-of-bino)
<http://rpubs.com/edwardcooper/titanic1>

[https://www.kaggle.com/mrisdal/exploring-survival-on-the-titanic/](https://www.kaggle.com/mrisdal/exploring-survival-on-the-titanic/report)
[report](https://www.kaggle.com/startupsci/titanic-data-science-solutions/comments)

[https://www.kaggle.com/startupsci/titanic-data-science-solutions/](https://www.kaggle.com/startupsci/titanic-data-science-solutions/comments)
[comments](https://www.newscientist.com/article/dn22119-sinking-the-titanic-women-and-children-f)

<https://www.newscientist.com/article/dn22119-sinking-the-titanic-women-and-children-f>

1. Introduction

The sinking of RMS Titanic brought to various machine learning competitions a quintessential dataset among others. After the “unsinkable” British passenger liner struck an iceberg in her maiden voyage on 15 April 1912 and was eventually wrecked, more than 1500 people perished. Decades of effort has been devoted to the study of the historic event, in which one major interest for statisticians is to predict possibility of survival given a number of characteristics, since there was clear account that some people (woman, children) were allowed to get on the lifeboat first.

There are several variants of Titanic data existed on the web, with primary source based on [Encyclopedia Titanica](#) (1999) founded by Philip Hind. This project is based on the most recent version with following columns available (table 1).

After appropriate formatting and cleaning, the data at hand recorded the survival status 2208 Titanic travelers alongside his/her gender, age, companions on board,

Table 1. Cleaned data with 2208 rows and 11 columns

Variable	Definition	Note
survived	Survival Status	0 = Lost, 1 = Saved
age	Age	In years, some infants had fractional values
gender	Gender	
class_dept	Class or Department	Passengers, Crew or Staff
nationality	Motherland	from wiki passenger list
title	Title	Extracted from name
spouse	# of spouse on board	
sibling	Number of siblings on board	
parent	Number of parents on board	
children	Number of children on board	

title, nationality, etc. There were 1496 victims and 712 survivors in total. Steps of data cleaning are elaborated in the [data](#) section in the appendix.

It is essential for every fruitful task of data analysis to first identify key questions of investigation that facilitates interpretation, however vague they are at the beginning. Then we can approach the core problem, filtering out trivialities, with statistical expression by abstraction. For our purposes, we could establish the following questions for which to quest

- To which degree is *Women and children first* policy respected? After the collision, the captain explicitly issued an order for women and children to be saved first.¹ Thus we should expect significantly higher proportion of females and children rescued than that in males and adults. If the opposite is true, that Titanic subjects behave more in line with the selfish *homo oeconomicus*, where everybody looked out for himself or herself and possibly even puts other people’s lives in danger, then people in their prime with physical superiority would see higher probability of survival. This requires us to study gender and age effect.
- Did socio-economic advantages mean better chance of survival? If this is the case, passengers with higher financial means, i.e. who live in the first class are more likely to survive. Similarly, passengers from second class will have a higher change of survival than third class people. Cabin class’s impact on survival status needs special notice here.
- For those who traveled alone with no companions (spouse, sibling, parent, children) on the vessel, is their survival possibility greater or less? On one hand, they are more likely to be in shortage of psychological and physical support. On the other hand, they would may be able to reach a life-saving decision faster without transaction cost and negotiation.
- Did English subjects receive any special care or given priority to aboard lifeboats? After all, Titanic was operated by British crew, and managed by British captain, masters and officers. Conversely, British nobility and elite
- Quantify interactions among various characteristics. Specifically, there are important interactions that need extra notice. For example, it has been widely studied in sociology and anthropology that human are sometimes driven by *procreation instinct* so that social norms would entail needs to protect females of reproductive age ([Frey, Savage, and Torgler 2009](#)).² Therefore, we could specify

¹Though there is no international maritime law enforcing this kind of chivalry.

²The average peak reproductive period in females is between the ages of 16 and 35.

and study the interaction between age and gender. Another typical interaction is between offspring and gender. *Parental investment* suggest that women on average invest more in caring for their offspring than males. In times of a disaster, higher opportunity cost will alert females with offspring more than others, and make them seek more aggressively for changes to secure the children as well as themselves.

This case study has been greatly inspired by Dr. Frank Harrell's similar example in his *Regression Modeling Strategies* (2015, Chapter 12) book, here I attempt to propose my understanding and interpretation of model development that is as original as possible. To ensure reproducibility, all the analysis is done in R (R Core Team 2020) and RStudio with code and text made public in this [repo](#). A brief summary of each section is listed below

- **Exploration.** Use descriptive statistics to examine data distribution characteristics, data missing patterns and relative effects, followed by redundancy analysis to study dependencies among predictors. Finish with nonparametric loess regression exploring nonlinear trends.
- **Modeling.** multiple imputation
- **Discussion.** Provide model-based explanation to answer former questions, combined with other research.
- **Conclusion.**

2. Exploration

2.1. Descriptive statistics and data processing

A graphical summary of the data is given by the `Hmisc::describe` function. For numerical variables, a inline histogram is produced alongside summary measures such as the number of missing values and the mean. For discrete variables, we focus on the number of categories and their relative frequency.

```
# print a summary for the data
t %>%
  describe() %>%
  latex(file = "", size = "small", center = "none")
```

11 Variables					2208 Observations									
survived														
n	missing	distinct	Info	Sum	Mean	Gmd								
2208	0	2	0.655	712	0.3225	0.4372								
age														
n	missing	distinct	Info	Mean	Gmd	.05	.10	.25	.50	.75	.90	.95		
1497	711	71	0.999	30.18	14.31	8	17	22	29	38	47	54		
lowest : 0.8 1.0 2.0 3.0 4.0, highest: 67.0 69.0 70.0 71.0 74.0														

gender											
n	missing	distinct									
2208	0	2									
Value	Female	Male									
Frequency	489	1719									
Proportion	0.221	0.779									
class_dept											
n	missing	distinct									
2208	0	5									
lowest :	1st	2nd	3rd	crew	staff	highest:	1st	2nd	3rd	crew	staff
Value	1st	2nd	3rd	crew	staff						
Frequency	321	270	709	822	86						
Proportion	0.145	0.122	0.321	0.372	0.039						
joined											
n	missing	distinct									
2208	0	4									
Value	Belfast	Cherbourg	Queenstown	Southampton							
Frequency	200	271	123	1614							
Proportion	0.091	0.123	0.056	0.731							
nationality											
n	missing	distinct									
2208	0	7									
lowest :	American	English	Finnish	Irish	Other	highest:	Finnish	Irish	Other	Swedish	Syrian
Value	American	English	Finnish	Irish	Other	Swedish	Syrian				
Frequency	246	1002	58	168	549	99	86				
Proportion	0.111	0.454	0.026	0.076	0.249	0.045	0.039				
title											
n	missing	distinct									
2208	0	4									
Value	Miss	Mr	Mrs	other							
Frequency	267	1590	212	139							
Proportion	0.121	0.720	0.096	0.063							
spouse											
n	missing	distinct	Info	Sum	Mean	Gmd					
2208	0	2	0.087	66	0.02989	0.05802					
sibling											
n	missing	distinct	Info	Mean	Gmd						
2208	0	4	0.138	0.05752	0.1103						
Value	0	1	2	3							
Frequency	2101	91	12	4							
Proportion	0.952	0.041	0.005	0.002							
parent											
n	missing	distinct	Info	Mean	Gmd						
2208	0	3	0.079	0.03804	0.07441						
Value	0	1	2								
Frequency	2148	36	24								
Proportion	0.973	0.016	0.011								
children											
n	missing	distinct	Info	Mean	Gmd						
2208	0	5	0.077	0.03895	0.07636						
lowest :	0	1	2	3	4	highest:	0	1	2	3	4
Value	0	1	2	3	4						
Frequency	2150	37	16	3	2						
Proportion	0.974	0.017	0.007	0.001	0.001						

There are several noteworthy patterns.³

³Though this may not be relevant to the model, it is still an surprising discovery that it wasn't until the late

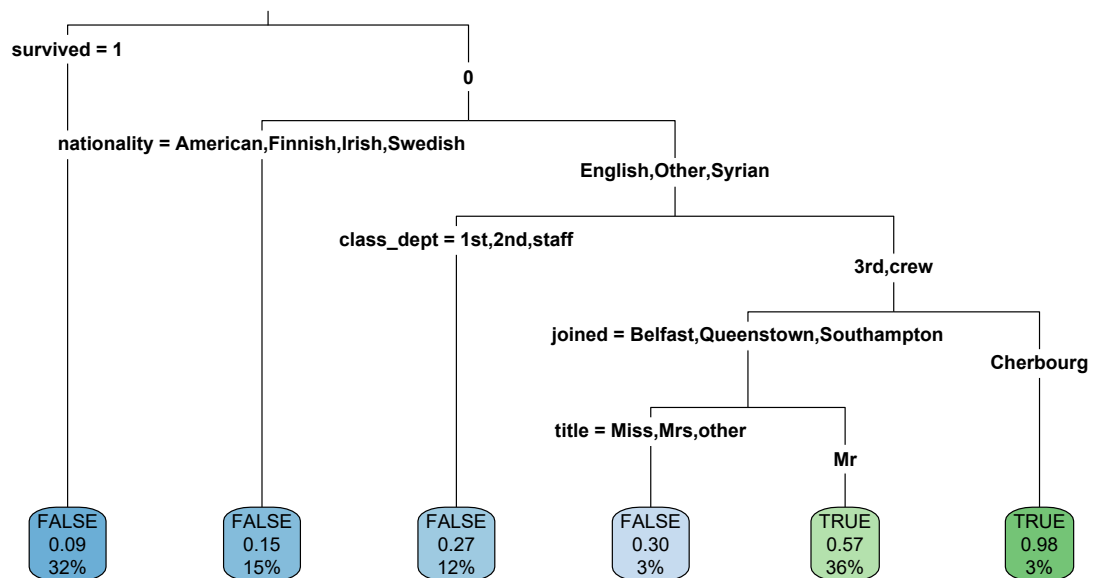


Figure 1. The decision tree for predicting `is.na(age)`, which finds strong patterns of missing related to class/department and gender (the Cherbourg node has very limited samples).

Of special importance is the `age` variable, which has roughly 30% missingness. On the other hand, it has a nice distribution with 80% known observations falling between 14 and 50. For further examination of patterns of missing data, we could fit a decision tree (figure 1) to predict which type of subject tend to have missing ages. Generally, for some third class male passenger or crew, age is mostly to miss.

```

na_tree <- rpart(factor(is.na(age)) ~ .,
  data = t %>% mutate(survived = as.factor(survived)) ,
  minbucket = 50)
rpart.plot::rpart.plot(na_tree, type = 3, cex = 0.6)

```

Back to other variables in descriptive statistics. Distributions of subject's companion on Titanic are all too narrow, as shown in figure 2. This motivates categorization since we will not lose too much information. Lastly, nearly half of the subjects are English. And if we focus on crew, the number rise to 85%.

Given this results, the final step in data munging is to dichotomize `spouse`, `parent`, `children` and `sibling` to denote if there is such relation. Thus we no longer have to deal with continuous predictors with poor distribution.

Univariate relationship between each independent variable and survival status is presented in figure 3. For each column, this is a anova-type plot with no control over confounding variables, though it may still assist us in determining how to spend degrees of freedom. If a predictor's effect on the response is strong, it's more likely that we need to spend more parameters on it. However, if a variable's effect appears to be weak, it could either result from a flat relationship with the response, or from

19th century that the idea of women traveling alone gained ground. As a result, there were nearly twice as many males passengers as females on Titanic. In fact, only 40% female passengers have no companion on the ship.



Figure 2. Few subjects have more than one companion in any of the 4 relations. Y axis on log scale.

nonlinearity and interaction among variables this plot fails to detect.

Finally, redundant analysis

Companion

2.2. *Loess regression for nonlinear pattern*

The loess is a common nonparametric regression method to study nonlinear relationship. In the case of binary response, the fitted value at $x = x_0$ is the proportion of positive cases near the neighborhood of x_0 ⁴. If the trend of a loess curve shows non-monotonicity, it is reasonable to include that nonlinearity relationship in the model, e.g., modeling the predictor with polynomial transformation or with splines.

Another important interaction, according to various literature, is related to cabin class (for passenger) and department (for crew and staff).

⁴

3. Modeling

the choice of model. In this setting, it is obvious that we would prefer probabilistic predictions to classification with output label 0 and 1, since we are placing emphasis upon the *tendency* of survival. And the true value of our model consist not in the decision on who will survive, but in what characteristics would increase or decrease the possibility of survival. The notion has ruled out most of the black box machine learning models for classification, say, random forest, support vector machines and neural network. Not only are they not intrinsically probability oriented, it is hard to

⁴with varying weights according to their distance to x_0

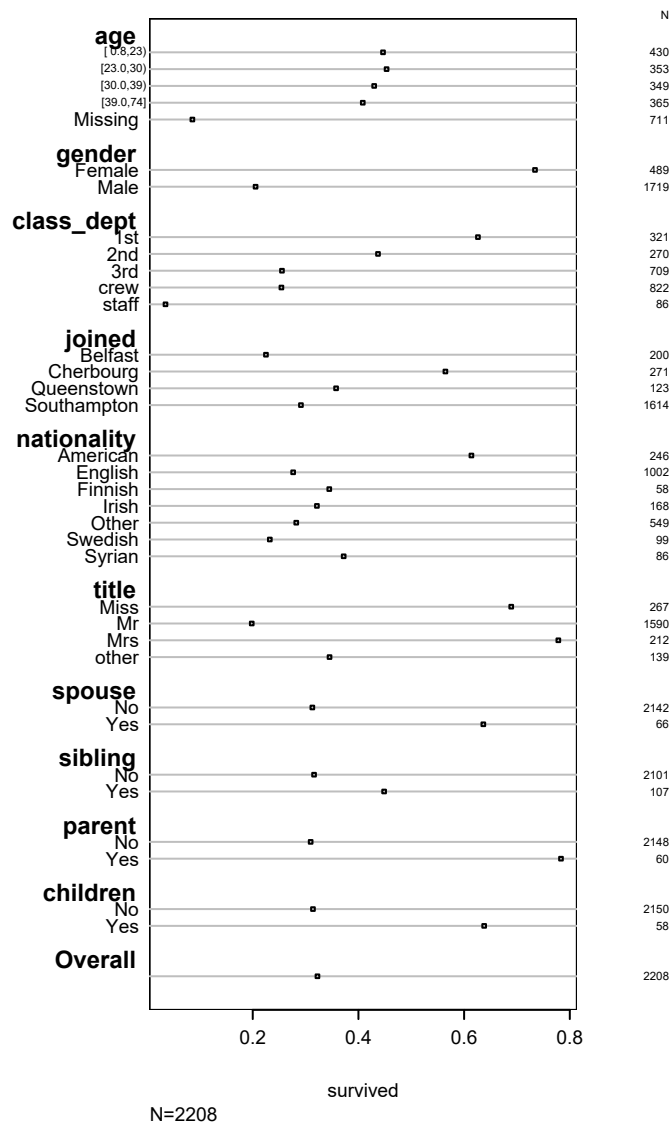


Figure 3. Summary of relationship between survival and each predictor

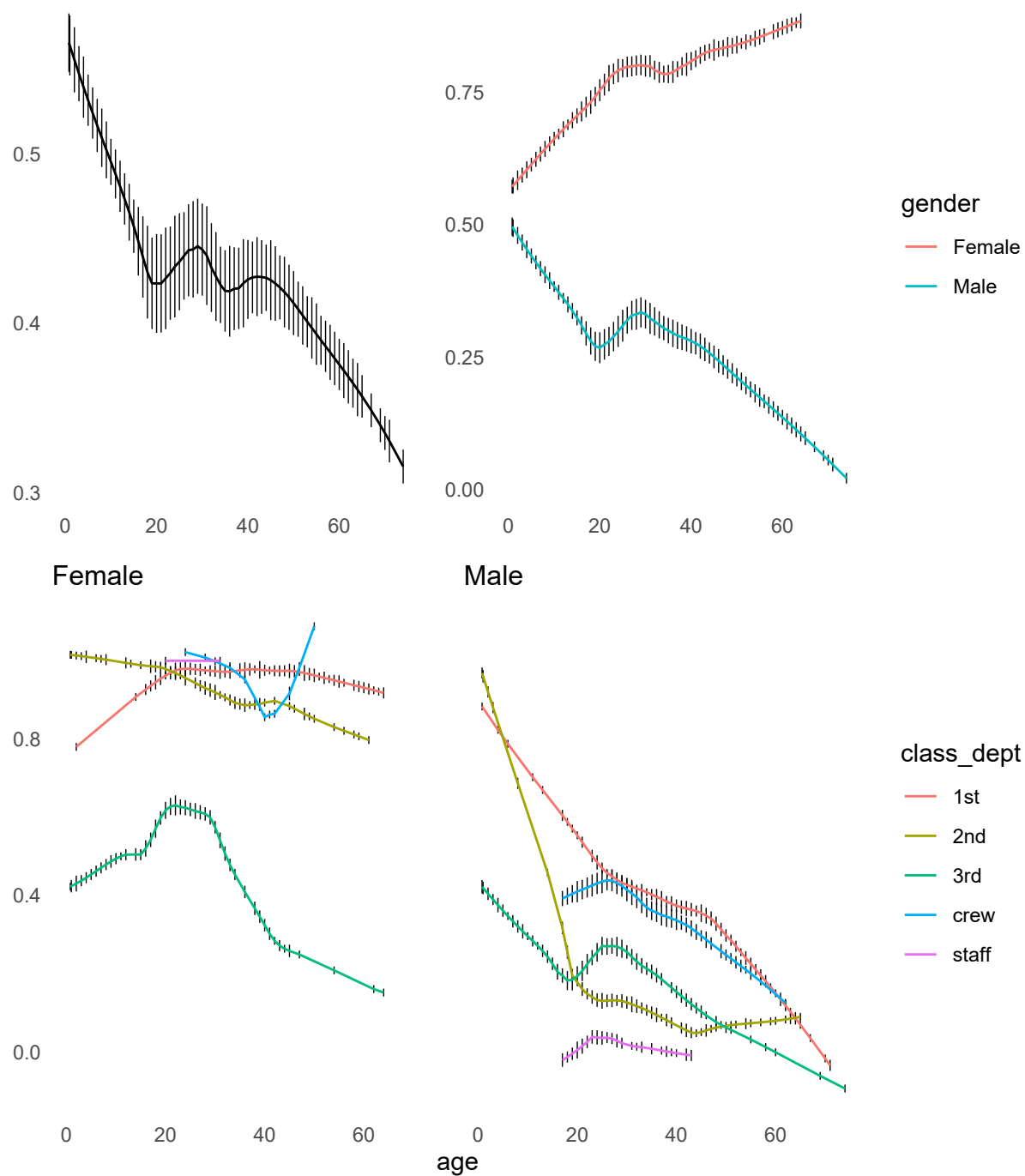


Figure 4. `loess` estimates of $P(\text{survived})$, with tick marks representing frequency counts within equal-width bins. Top left panel shows the nonlinear relationship between age and survival status without controlling confounding variables. Other plots give estimates under stratification by sex and class/department.

interpret main effects and interactions as everything seems to be interacted with one another.

3.1. Saturated model

First and foremost,

The limiting sample size for binary outcome would be the number of minority class, in our case 712. Using the 15:1 rule, that will give us some confidence spending roughly 47 parameters or degrees of freedom.

This plot is used to identify possibly flat relationship between predictor and response. While misuse of this plot would be checking nonlinearity. Even with spline transformation and large corrected χ^2 there is no guarantee for nonlinearity.

In this sense, the saturated model could provide rough guidance

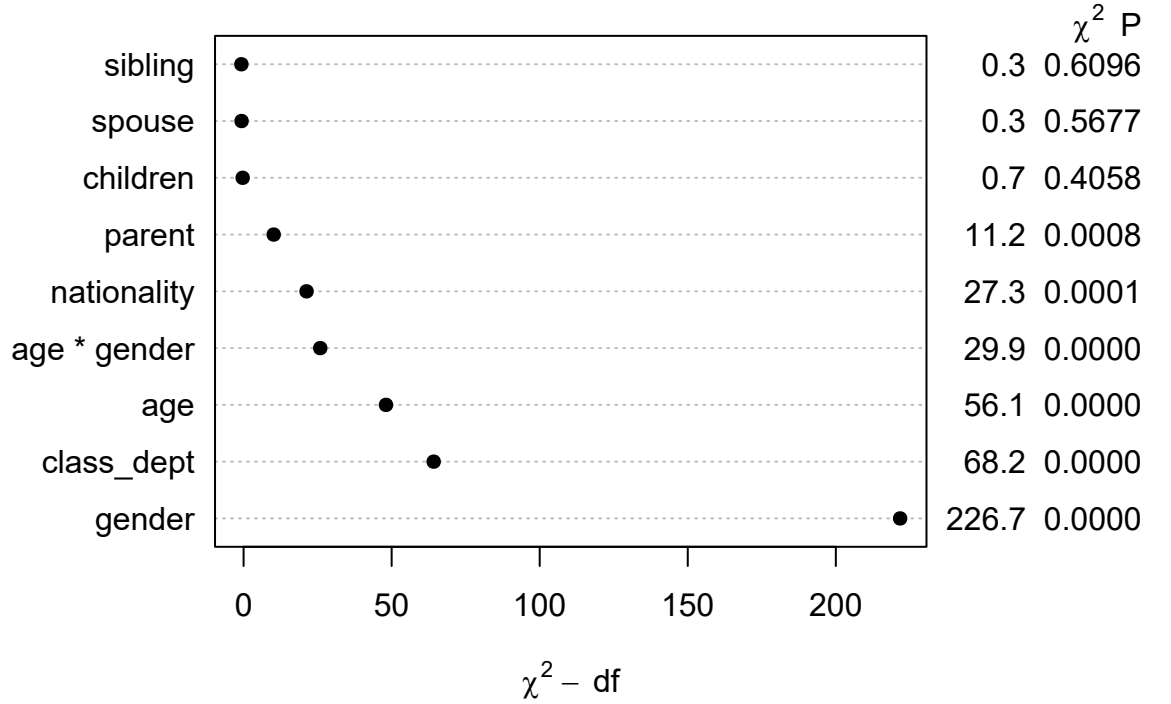
Table 2. d.f. budget in the saturated model

age	sibling	parent	children
4	1	1	1

hypothesis testing
anova plot

Table 3. Wald Statistics for `survived`

	χ^2	d.f.	<i>P</i>
age (Factor+Higher Order Factors)	56.09	8	<0.0001
<i>All Interactions</i>	29.89	4	<0.0001
<i>Nonlinear (Factor+Higher Order Factors)</i>	20.81	6	0.0020
gender (Factor+Higher Order Factors)	226.72	5	<0.0001
<i>All Interactions</i>	29.89	4	<0.0001
class_dept	68.20	4	<0.0001
nationality	27.27	6	0.0001
spouse	0.33	1	0.5677
sibling	0.26	1	0.6096
parent	11.17	1	0.0008
children	0.69	1	0.4058
age \times gender (Factor+Higher Order Factors)	29.89	4	<0.0001
<i>Nonlinear</i>	12.60	3	0.0056
<i>Nonlinear Interaction : f(A,B) vs. AB</i>	12.60	3	0.0056
TOTAL NONLINEAR	20.81	6	0.0020
TOTAL NONLINEAR + INTERACTION	40.05	7	<0.0001
TOTAL	305.90	23	<0.0001



3.2. Multiple imputation

The pooled estimates are obtained by averaging over m fitted model based on one piece of multiple imputation. The variance-covariance matrix T is calculated using Rudin's rule

$$T = \frac{1}{m} \sum_{i=1}^m U_i + \left(1 + \frac{1}{m}\right) B$$

where U_i is the estimated complete-data variance-covariance matrix in each imputation, and B the estimated variance-covariance matrix between the m complete-data estimates. Here we see the one major advantage of multiple imputation over single imputation is that not only does its variance estimates accounts for sampling variability, but also for the extra variance caused by missing values and finite number of imputations.

There are some simple workarounds

- complete-case analysis: That is, we delete all incomplete observations. Needless to say this will translate into a major harm on sample size since over 60% of **boat** are missing, not to mention other columns. Even if we remove **boat** and then delete rows with missing **age** we still lose over 1/5 of data. Moreover, figures in 2 have shed light on the relatively strong influence of **age** on survival. Also,

the deletion of incomplete observations assumes data are missing completely at random (MCAR). When it's not the case, this could severely bias estimates of coefficients (Van Buuren 2018)

- single imputation:
- multiple imputation

3.3.

3.4. *Model fitting and penalization*

<https://www.encyclopedia-titanica.org/community/threads/passengers-who-spoke-other-languages.20103/>

Since the crew's instructions (in English) tended to be along the lines of "Wait down here for further orders" a lack of understanding might well have saved many lives. Also many of the immigrants in 3rd Class were traveling in family or neighbourhood groups which included at least one English-speaker (often an established immigrant returning to the US from a visit back home) who could act as their spokesperson.

3.5. *Validation and calibration*

There will not be another Titanic, and any model on Titanic will not be used for prediction. Therefore, the goal of model validation is primarily to provide quantify the degree of overfitting with various bias-corrected measures.

In the award-winning solution to this legendary dataset presented by IBM Watson, they used a holdout sample to validate their model. <https://www.fharrell.com/post/split-val/>

As a integral component of model validation, calibration aims to gauge the concordance between predicted values and observed data.

3.6. *Diagnostics and interpretation*

influence

`which.influence`

4. Discussion

The most decisive explanation for such effect is that first-class passengers had better access to information about the imminent danger and were aware that the lifeboats were located close to the first class cabins. Thus, their marginal effort costs to survive were lower. In contrast, most third-class passengers had no idea where the lifeboats were located (safety drills for all passengers were introduced after the Titanic disaster), and they did not know how to reach the upper decks where the lifeboats were stowed.

Wyn Craig Wade: there was a class culture on Titanic akin to the notion of a "culture of poverty

Undoubtedly, the worst barriers were the ones within the steerage passengers themselves. Years of conditioning as third-class citizens led a great many of them to give up hope as soon as the crisis became evident ... Barriers to steerage? Yes, but of a kind less indictable to the White Star Line than to the whole of civilization.

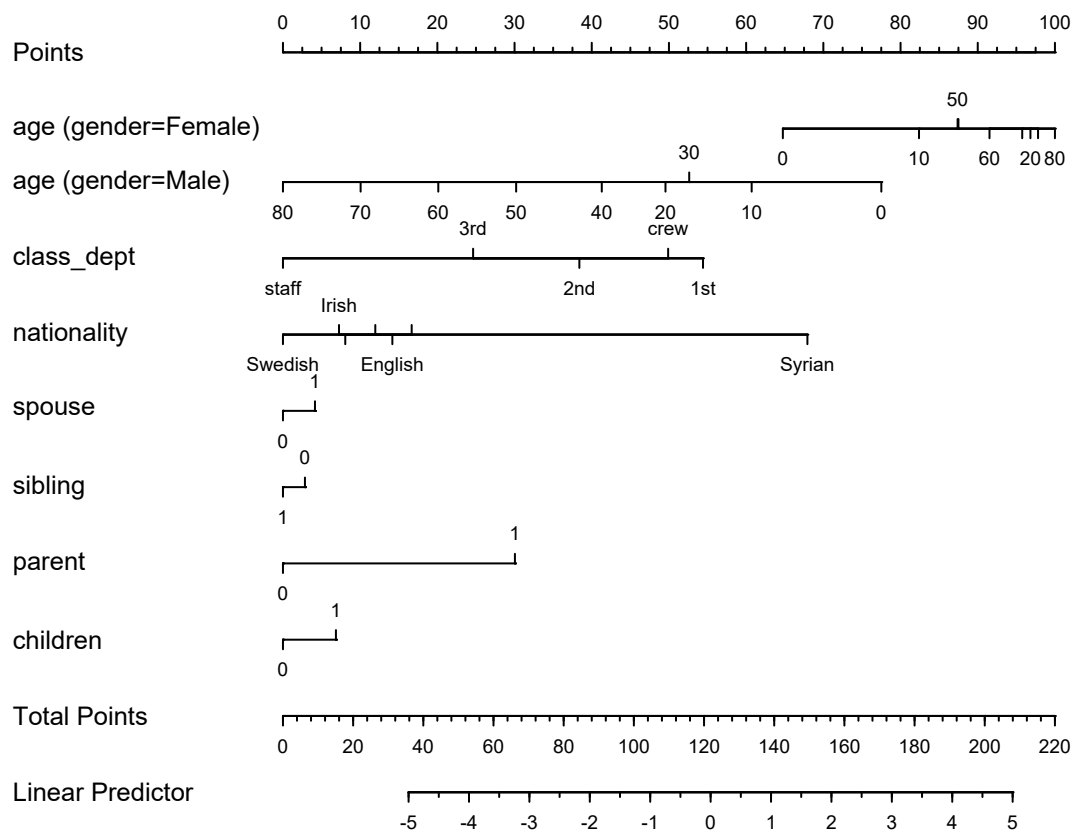


Figure 5. nomogram

A more detailed explanation of some of these measures is presented in the [appendix](#).

Women and children first only for higher class passengers. If you are a third class female

5. Conclusion

Appendix A. Data

The source data is accessed on [Encyclopedia Titanica](#), a leading archive on titanic facts. In contrast to the famous titanic dataset (known as `titanic3`) distributed by [kaggle](#) for introductory level machine learning practices, the case study uses a more up-to-date and complete dataset in the following ways

- **Larger sample size.** Our data includes crew and staff members alongside passengers, while `titanic3` only incorporate passenger information. We do not use a separate test set approach for validation either. As a result, the sample size is about 2.5 times larger.
- **More columns.** Additional variables such as role on the ship, nationality and occupation are added. A major difference is made by separating the travel companion data into four distinct columns: number of parents, children, sibling and spouses that each passenger traveled with. These were combined into two columns before.
- **More accurate.** `titanic3` was an effort to study Titanic in the 20th century, lastly updated and improved by Thomas Cason in 1999. The data has been constantly revised, many errors corrected, many missing ages filled in, and new variables created. Now it reflects the state of the data as of 21 October 2020.

The data cleaning process involves using appropriate data types, creating new features, adjusting levels for categorical variable and excluding irrelevant columns. Code can be found at [clean.R](#).

`title` is extracted through each person’s name with regular expressions and then collapsed into 4 levels.⁵

Passengers are classified according to their cabin class. Others on the vessel fall into one of crew and staff members. Crew includes victualling crew⁶, engineering crew, deck crew and officers, substitute crew and guarantee group. Staff members include restaurant staff and orchestra.

Rare nationality (lower than 50 people) is collapsed.

Age information is presented as non-missing on the surface yet there is an indicator column representing when a person’s age is only approximate and cannot be fully determined from current facts. These inaccurate age have been assigned NA. There were also ten subjects whose four companion variables were all explicitly missing. For simplicity, the mode 0 is filled in. Therefore, the problem of missing data is reduced to univariate missing of `age`.

Variables we do not utilize in this project includes name, date of birth and death, lifeboat number⁷, fare, and cabin number.⁸

⁵For example, the title for passenger “Abbing, Mr Anthony” is “Mr”.

⁶crew in charge of food, housekeeping, laundry, room service, etc.

⁷There were 9 recorded passengers who got on the lifeboat yet died before reaching Carpathia, another RMS which spearheaded the rescue of Titanic survivors. There were also 13 passengers who survived with no boat information documented, and this is most likely due to data quality issues after looking up on Encyclopedia Titanica. Even with these exceptions, whether a passenger got on a lifeboat yields perfect prediction on his/her survival. If one fits a logistic regression model on survival based on whether `boat` is missing, the apparent accuracy will be nearly 1. In this sense `boat` is more the result of survival, rather than a cause.

⁸While some study used this attribute to find cabin locations, its large amount of missingness could be a major source of complexity.

Appendix B. Model formula

The formula for our binary logistic model

Appendix C. Measures used in valiation

Somer's D_{xy} index is a calibration measure, which is the rank correlation between predicted and actual response. It has a close relationship with the C index

$$D_{xy} = 2(c - 0.5)$$

Appendix D. Original Computing Environment

```
sessionInfo()
```

```
R version 4.0.2 (2020-06-22)
```

```
Platform: x86_64-w64-mingw32/x64 (64-bit)
```

```
Running under: Windows 10 x64 (build 18362)
```

```
Matrix products: default
```

```
locale:
```

```
[1] LC_COLLATE=English_United States.1252
```

```
[2] LC_CTYPE=English_United States.1252
```

```
[3] LC_MONETARY=English_United States.1252
```

```
[4] LC_NUMERIC=C
```

```
[5] LC_TIME=English_United States.1252
```

```
system code page: 936
```

```
attached base packages:
```

```
[1] stats      graphics  grDevices  utils      datasets  methods    base
```

```
other attached packages:
```

```
[1] rpart_4.1-15      patchwork_1.0.1  mice_3.11.0      rms_6.0-1
```

```
[5] SparseM_1.78      Hmisc_4.4-1      Formula_1.2-4    survival_3.1-12
```

```
[9] lattice_0.20-41  ggplot2_3.3.2    dplyr_1.0.2
```

```
loaded via a namespace (and not attached):
```

```
[1] tidyr_1.1.2      splines_4.0.2    assertthat_0.2.1
```

```
[4] latticeExtra_0.6-29  yaml_2.2.1      pillar_1.4.6
```

```
[7] backports_1.1.10    quantreg_5.74    glue_1.4.2
```

```
[10] digest_0.6.26      RColorBrewer_1.1-2  checkmate_2.0.0
```

```
[13] colorspace_1.4-1    sandwich_3.0-0    htmltools_0.5.0
```

```
[16] Matrix_1.2-18      conquer_1.0.2     pkgconfig_2.0.3
```

```
[19] broom_0.7.2        bookdown_0.21     purrr_0.3.4
```

```
[22] mvtnorm_1.1-1      scales_1.1.1      jpeg_0.1-8.1
```

[25]	MatrixModels_0.4-1	htmlTable_2.1.0	tibble_3.0.4
[28]	rticles_0.16.1	farver_2.0.3	generics_0.0.2
[31]	ellipsis_0.3.1	TH.data_1.0-10	withr_2.3.0
[34]	nnet_7.3-14	cli_2.1.0	magrittr_1.5
[37]	crayon_1.3.4	polyspline_1.1.19	evaluate_0.14
[40]	fansi_0.4.1	nlme_3.1-148	MASS_7.3-51.6
[43]	foreign_0.8-80	tools_4.0.2	data.table_1.13.2
[46]	hms_0.5.3	lifecycle_0.2.0	matrixStats_0.57.0
[49]	multcomp_1.4-14	stringr_1.4.0	rpart.plot_3.0.9
[52]	munSELL_0.5.0	cluster_2.1.0	compiler_4.0.2
[55]	rlang_0.4.8	grid_4.0.2	rstudioapi_0.11
[58]	htmlwidgets_1.5.2	base64enc_0.1-3	labeling_0.4.2
[61]	rmarkdown_2.5	gtable_0.3.0	codetools_0.2-16
[64]	R6_2.4.1	gridExtra_2.3	zoo_1.8-8
[67]	knitr_1.30	readr_1.4.0	stringi_1.5.3
[70]	Rcpp_1.0.5	vctrs_0.3.4	png_0.1-7
[73]	tidyselect_1.1.0	xfun_0.18	

References

- Allaire, JJ, Yihui Xie, Jonathan McPherson, Javier Luraschi, Kevin Ushey, Aron Atkins, Hadley Wickham, Joe Cheng, Winston Chang, and Richard Iannone. 2020a. *rmarkdown: Dynamic Documents for R*. R package version 2.5, <https://github.com/rstudio/rmarkdown>.
- Allaire, JJ, Yihui Xie, R Foundation, Hadley Wickham, Journal of Statistical Software, Ramnath Vaidyanathan, Association for Computing Machinery, et al. 2020b. *rticles: Article Formats for R Markdown*. R package version 0.16.1, <https://github.com/rstudio/rticles>.
- Frey, Bruno S, David A Savage, and Benno Torgler. 2009. "Surviving the Titanic disaster: economic, natural and social determinants." .
- Harrell, Frank E, Jr. 2020. *Hmisc: Harrell Miscellaneous*. R package version 4.4-1, <https://CRAN.R-project.org/package=Hmisc>.
- Harrell Jr, Frank E. 2015. *Regression modeling strategies: with applications to linear models, logistic and ordinal regression, and survival analysis*. Springer.
- Harrell, Jr., Frank E. 2020. *rms: Regression Modeling Strategies*. R package version 6.0-1, <https://CRAN.R-project.org/package=rms>.
- Hind, Philip. 1999. <https://www.encyclopedia-titanica.org/>.
- Koenker, Roger, and Pin Ng. 2019. *SparseM: Sparse Linear Algebra*. R package version 1.78, <http://www.econ.uiuc.edu/~roger/research/sparse/sparse.html>.
- Milborrow, Stephen. 2020. *rpart.plot: Plot rpart Models: An Enhanced Version of plot.rpart*. R package version 3.0.9, <http://www.milbo.org/rpart-plot/index.html>.
- R Core Team. 2020. *R: A Language and Environment for Statistical Computing*. Vienna, Austria: R Foundation for Statistical Computing. <https://www.R-project.org/>.
- Roecker, Ellen B. 1991. "Prediction error and its estimation for subset-selected models." *Technometrics* 33 (4): 459–468.
- Sarkar, Deepayan. 2008. *Lattice: Multivariate Data Visualization with R*. New York: Springer. ISBN 978-0-387-75968-5, <http://lmdvr.r-forge.r-project.org>.
- Sarkar, Deepayan. 2020. *lattice: Trellis Graphics for R*. R package version 0.20-41, <http://lattice.r-forge.r-project.org/>.
- Terry M. Therneau, and Patricia M. Grambsch. 2000. *Modeling Survival Data: Extending the Cox Model*. New York: Springer.
- Therneau, Terry, and Beth Atkinson. 2019. *rpart: Recursive Partitioning and Regression Trees*. R package version 4.1-15, <https://CRAN.R-project.org/package=rpart>.

- Therneau, Terry M. 2020. *survival: Survival Analysis*. R package version 3.1-12, <https://github.com/therneau/survival>.
- Van Buuren, Stef. 2018. *Flexible imputation of missing data*. CRC press.
- van Buuren, Stef, and Karin Groothuis-Oudshoorn. 2011. “mice: Multivariate Imputation by Chained Equations in R.” *Journal of Statistical Software* 45 (3): 1–67. <https://www.jstatsoft.org/v45/i03/>.
- van Buuren, Stef, and Karin Groothuis-Oudshoorn. 2020. *mice: Multivariate Imputation by Chained Equations*. R package version 3.11.0, <https://CRAN.R-project.org/package=mice>.
- Wickham, Hadley. 2016. *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York. <https://ggplot2.tidyverse.org>.
- Wickham, Hadley, Winston Chang, Lionel Henry, Thomas Lin Pedersen, Kohske Takahashi, Claus Wilke, Kara Woo, Hiroaki Yutani, and Dewey Dunnington. 2020a. *ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics*. R package version 3.3.2, <https://CRAN.R-project.org/package=ggplot2>.
- Wickham, Hadley, Romain François, Lionel Henry, and Kirill Müller. 2020b. *dplyr: A Grammar of Data Manipulation*. R package version 1.0.2, <https://CRAN.R-project.org/package=dplyr>.
- Xie, Yihui. 2014. “knitr: A Comprehensive Tool for Reproducible Research in R.” In *Implementing Reproducible Computational Research*, edited by Victoria Stodden, Friedrich Leisch, and Roger D. Peng. Chapman and Hall/CRC. ISBN 978-1466561595, <http://www.crcpress.com/product/isbn/9781466561595>.
- Xie, Yihui. 2015. *Dynamic Documents with R and knitr*. 2nd ed. Boca Raton, Florida: Chapman and Hall/CRC. ISBN 978-1498716963, <https://yihui.org/knitr/>.
- Xie, Yihui. 2016. *bookdown: Authoring Books and Technical Documents with R Markdown*. Boca Raton, Florida: Chapman and Hall/CRC. ISBN 978-1138700109, <https://github.com/rstudio/bookdown>.
- Xie, Yihui. 2020a. *bookdown: Authoring Books and Technical Documents with R Markdown*. R package version 0.21, <https://github.com/rstudio/bookdown>.
- Xie, Yihui. 2020b. *knitr: A General-Purpose Package for Dynamic Report Generation in R*. R package version 1.30, <https://yihui.org/knitr/>.
- Xie, Yihui, J.J. Allaire, and Garrett Golemund. 2018. *R Markdown: The Definitive Guide*. Boca Raton, Florida: Chapman and Hall/CRC. ISBN 9781138359338, <https://bookdown.org/yihui/rmarkdown>.
- Xie, Yihui, Christophe Dervieux, and Emily Riederer. 2020. *R Markdown Cookbook*. Boca Raton, Florida: Chapman and Hall/CRC. ISBN 9780367563837, <https://bookdown.org/yihui/rmarkdown-cookbook>.
- Zeileis, Achim, and Yves Croissant. 2010. “Extended Model Formulas in R: Multiple Parts and Multiple Responses.” *Journal of Statistical Software* 34 (1): 1–13.
- Zeileis, Achim, and Yves Croissant. 2020. *Formula: Extended Model Formulas*. R package version 1.2-4, <https://CRAN.R-project.org/package=Formula>.