

An exploration of data: prediction, modeling and displays

Karsten Maurer

**Department of Statistics
Miami University
Oxford, OH 45056
maurerkt@miamioh.edu**

URL: kmaurer.github.io

Department Social Media

Twitter: @STADeprMiamiOH @statsandstories

Podcast: www.statsandstories.net

A black and white photograph of the RMS Titanic, showing its four funnels and the ship's hull. The ship is viewed from a side-on perspective, moving towards the right. The funnels are dark with light-colored tops. The hull is dark, and the upper decks are visible. The ship is on a calm sea under a light sky.

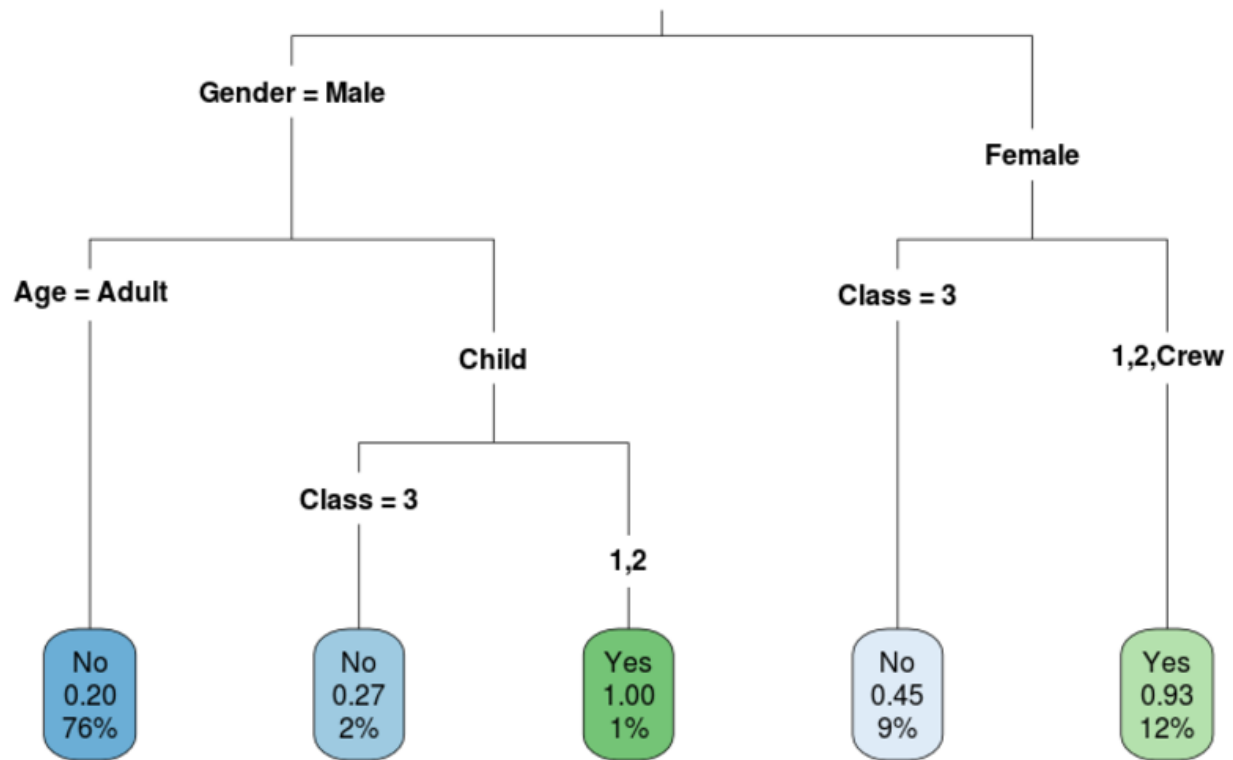
Data: In the paper, The "Unusual Episode" Data Revisited published in the *Journal of Statistics Education* vol.3, no.3 (1995), records for 2201 passengers and crew were recorded with their ticket status (the Class variable), Age (categorized as Adult/Child), Gender (Female/Male) and whether they survived the sinking. 15 of the 2201 passengers/crew were randomly removed from the record and summary tables of the remaining 2186 passengers/crew is included below.

1. Which variables appear to influence a person's survival?
2. On the next page is a list of the 15 people removed from the record. Your goal is to:
 - A) Predict whether each of the 15 people survived.
 - B) Assign a probability/percentage on the likelihood they survived.

2

<u>Person</u>	<u>Class</u>	<u>Age</u>	<u>Gender</u>
219	1st	Adult	Female
566	2nd	Adult	Female
602	2nd	Child	Female
633	3rd	Adult	Male
815	3rd	Adult	Male
866	3rd	Adult	Male
1104	3rd	Adult	Female
1122	3rd	Adult	Female
1402	Crew	Adult	Male
1407	Crew	Adult	Male
1672	Crew	Adult	Male
1854	Crew	Adult	Male
2025	Crew	Adult	Male
2097	Crew	Adult	Male
2135	Crew	Adult	Male

Suppose we built a statistical model ... a classification tree was produced below based on a training set of 2186 passengers. (STA 333, STA 467)



This could be applied to the test set of 15 passengers that were sampled.

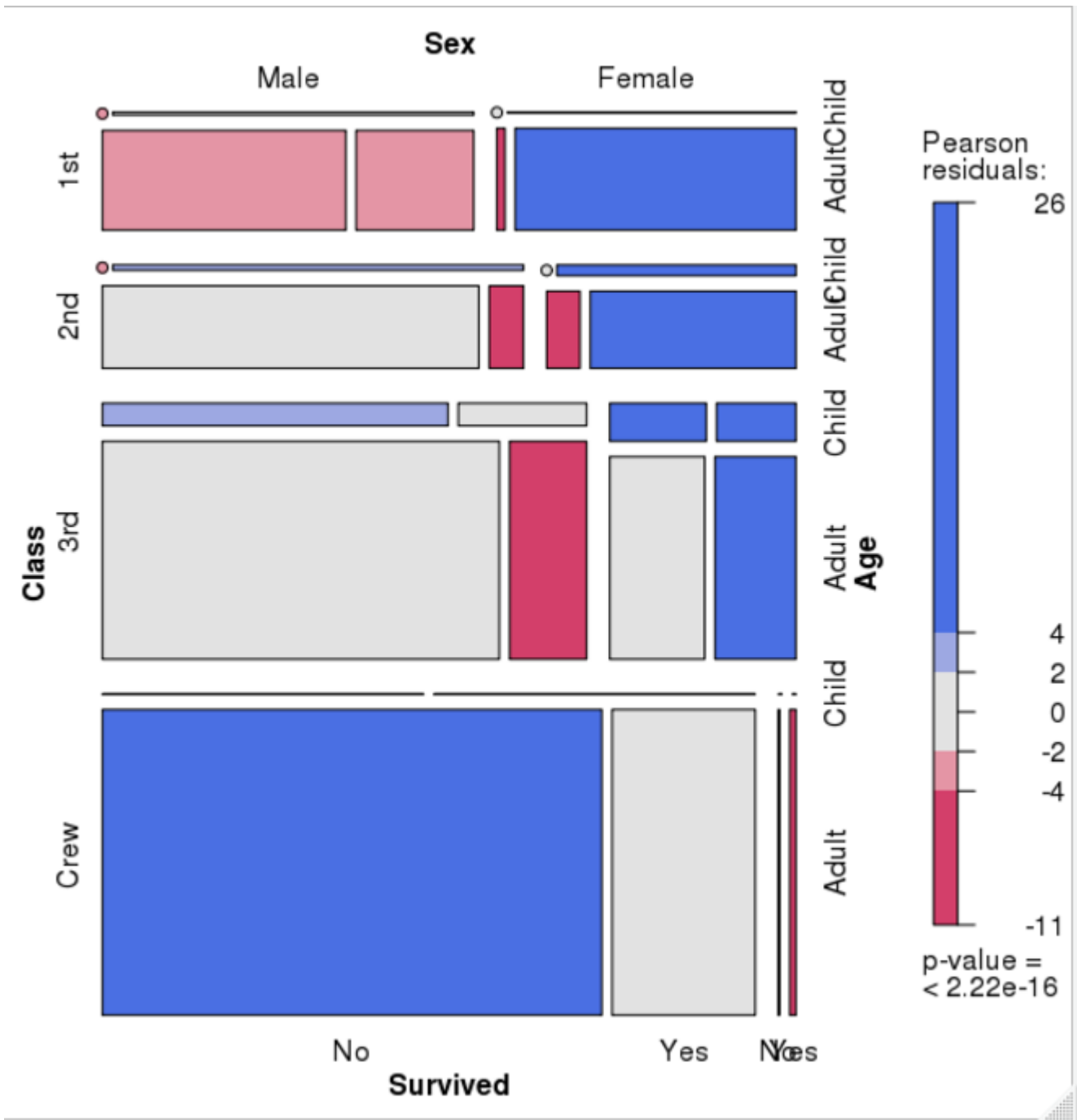
Predictions based on classification tree

<u>Person</u>	<u>Class</u>	<u>Age</u>	<u>Gender</u>	<u>P(Not)</u>	<u>P(Survived)</u>	<u>Predict</u>
219	1st	Adult	Female	0.07380074	0.9261993	Yes
566	2nd	Adult	Female	0.07380074	0.9261993	Yes
602	2nd	Child	Female	0.07380074	0.9261993	Yes
633	3rd	Adult	Male	0.79782740	0.2021726	No
815	3rd	Adult	Male	0.79782740	0.2021726	No
866	3rd	Adult	Male	0.79782740	0.2021726	No
1104	3rd	Adult	Female	0.54639175	0.4536082	No
1122	3rd	Adult	Female	0.54639175	0.4536082	No
1402	Crew	Adult	Male	0.79782740	0.2021726	No
1407	Crew	Adult	Male	0.79782740	0.2021726	No
1672	Crew	Adult	Male	0.79782740	0.2021726	No
1854	Crew	Adult	Male	0.79782740	0.2021726	No
2025	Crew	Adult	Male	0.79782740	0.2021726	No
2097	Crew	Adult	Male	0.79782740	0.2021726	No
2135	Crew	Adult	Male	0.79782740	0.2021726	No

<u>Person</u>	<u>Class</u>	<u>Age</u>	<u>Gender</u>	<u>P (Not)</u>	<u>P (Survived)</u>	<u>Prediction</u>	<u>Truth</u>
219	1st	Adult	Female	0.07380074	0.9261993	Yes	Yes
566	2nd	Adult	Female	0.07380074	0.9261993	Yes	Yes
602	2nd	Child	Female	0.07380074	0.9261993	Yes	Yes
633	3rd	Adult	Male	0.79782740	0.2021726	No	Yes
815	3rd	Adult	Male	0.79782740	0.2021726	No	No
866	3rd	Adult	Male	0.79782740	0.2021726	No	No
1104	3rd	Adult	Female	0.54639175	0.4536082	No	Yes
1122	3rd	Adult	Female	0.54639175	0.4536082	No	Yes
1402	Crew	Adult	Male	0.79782740	0.2021726	No	Yes
1407	Crew	Adult	Male	0.79782740	0.2021726	No	Yes
1672	Crew	Adult	Male	0.79782740	0.2021726	No	No
1854	Crew	Adult	Male	0.79782740	0.2021726	No	No
2025	Crew	Adult	Male	0.79782740	0.2021726	No	No
2097	Crew	Adult	Male	0.79782740	0.2021726	No	No
2135	Crew	Adult	Male	0.79782740	0.2021726	No	No

Our predictions of the 15 passengers that were sampled wasn't perfect (10 of 15 classified correctly).

Tables can be tough to process. Can we visualize this? (STA 404)



Visual cues in this Mosaic Plot?

- Size of boxes
- Color of boxes

Visualizing data ...

The Joy of Statistics – Hans Rosling

<https://www.youtube.com/watch?v=jbkSRLYSojo>

As you watch this video, please record the following information:

What variables were presented?

What graphical characteristic (aesthetic trait) was mapped to each variable?

Rstudio.miamioh.edu -> exploring gapminder data

Data-visualization-exploration-DEMO-04oct17.R

Studying at Miami University

- Math & Stat Degrees ([B.S. Math & Stat](#), [B.S. Stat](#))
 - Foundation in mathematics
 - Statistical modeling
 - Data handling and visualization
- [Analytics Co-Major](#)
 - Complements the B.S. Math & Statistics & B.S. Statistics very well
- [Actuarial Science Minor](#) (actuarial science club: Dr. Miljkovic)
 - Complements B.S. degrees and satisfies related hours & thematic sequences
- [Miami University StatHawks](#) (partners with Pi Mu Epsilon for some activities)
 - Student Chapter of the American Statistical Association
 - Can join on the Hub - Events throughout the fall (movie night, trivia, speakers)
- [Center for Analytics and Data Science](#) (CADS)
 - DataFest - weekend of April 6-8, 2018

THANK YOU!

Questions?