

FactSheet:: Titanic Data

M	Topic & Assignment																								
M2	Titanic data mining analysis																								
A.	Background and overviews <ul style="list-style-type: none">• https://www.rdocumentation.org/packages/titanic/versions/0.1.0• https://www.kaggle.com/competitions/titanic/overview• https://www.encyclopedia-titanica.org/ <p>The Titanic DataFrames describe the survival status of individual Titanic passengers, not the crew, with ages for ~half the passengers. One of the original sources is Eaton & Haas (1994) Titanic: Triumph and Tragedy, Patrick Stephens Ltd includes a passenger list created by many researchers and edited by Michael A. Findlay [1].</p>																								
B.	Interesting models - built in R code for display convenience <pre>> data <- read.csv('titanic.csv')</pre> <ul style="list-style-type: none">• # Linear regression model<pre>• model <- lm(survived ~ age + sex + pclass + sibsp + parch, data = data)</pre>• Binomial Predicting survival based on age, sex, and passenger class<pre>• model <- glm(survived ~ age + sex + pclass, data = titanic, family = binomial)</pre>• Poisson - Predicting the count of siblings/spouses based on passenger age<pre>• model <- glm(sibsp ~ age, data = titanic, family = poisson) summary(model)</pre>• Neg.Binomial - Predict count of parents/children by passenger age and sex<pre>• model <- glm.nb(parch ~ age + sex, data = titanic) summary(model)</pre>																								
C.	Data <class.github> <ul style="list-style-type: none">• raw data; unsplit and preprocessed [source: https://hbiostat.org/data/ <titanic.3>• train, test; from kaggle																								
D.	Data dictionary <table><tr><td>passengerid</td><td>sequential unique id</td></tr><tr><td>survived</td><td>0=no, 1=yes</td></tr><tr><td>pclass</td><td>1,2,3:passenger class (1st, 2nd, 3rd); proxy for socio-economic class</td></tr><tr><td>name</td><td>Christian name</td></tr><tr><td>sex</td><td>male, female</td></tr><tr><td>age</td><td>00, NA, blank. in years; some infants w fractional values</td></tr><tr><td>sibsp</td><td>number of siblings and spouses aboard</td></tr><tr><td>parch</td><td><parent.child> #parents or chil</td></tr><tr><td>ticket</td><td>alpha, numeric, character</td></tr><tr><td>fare</td><td>0.0000 decimals</td></tr><tr><td>cabin</td><td>C#, blank,</td></tr><tr><td>embarked</td><td>C, Q, S <Cherbourg, Southampton, and Queenstown></td></tr></table> <p>References:</p> <ol style="list-style-type: none">1. Harrell Jr, F.E.,(2002). Titanic data, Vanderbuilt biostatistics datasets. Vanderbilt University. Retrieved from: https://hbiostat.org/data/repo/titanic.html. Retrieved on 05.15.2023.	passengerid	sequential unique id	survived	0=no, 1=yes	pclass	1,2,3:passenger class (1st, 2nd, 3rd); proxy for socio-economic class	name	Christian name	sex	male, female	age	00, NA, blank. in years; some infants w fractional values	sibsp	number of siblings and spouses aboard	parch	<parent.child> #parents or chil	ticket	alpha, numeric, character	fare	0.0000 decimals	cabin	C#, blank,	embarked	C, Q, S <Cherbourg, Southampton, and Queenstown>
passengerid	sequential unique id																								
survived	0=no, 1=yes																								
pclass	1,2,3:passenger class (1st, 2nd, 3rd); proxy for socio-economic class																								
name	Christian name																								
sex	male, female																								
age	00, NA, blank. in years; some infants w fractional values																								
sibsp	number of siblings and spouses aboard																								
parch	<parent.child> #parents or chil																								
ticket	alpha, numeric, character																								
fare	0.0000 decimals																								
cabin	C#, blank,																								
embarked	C, Q, S <Cherbourg, Southampton, and Queenstown>																								

