

Instructions:

- => If course materials brought you here, scroll or use links to sections.
- => Use links on downloaded.pdf! In git, download arrow is on the right above doc visual.

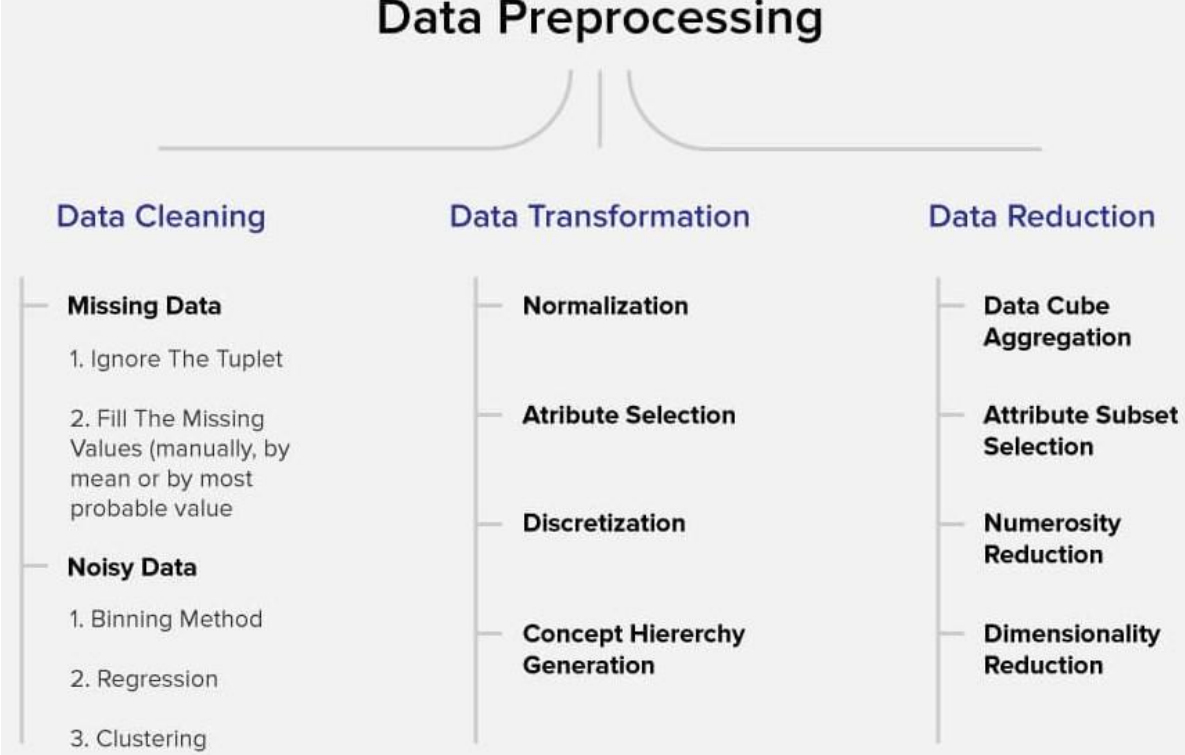
<a href="#">Titanic Data</a>	<a href="#">PreProcess + Supervised &amp; Unsupervised</a>		
------------------------------	--	--	--

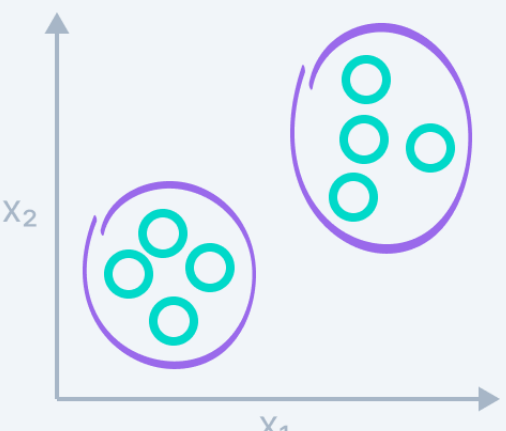
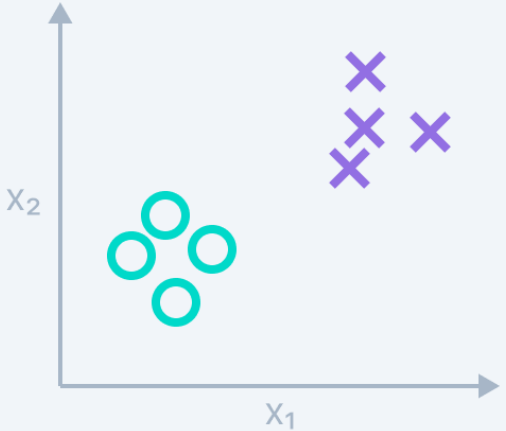
M.2.Assignment - Titanic Data

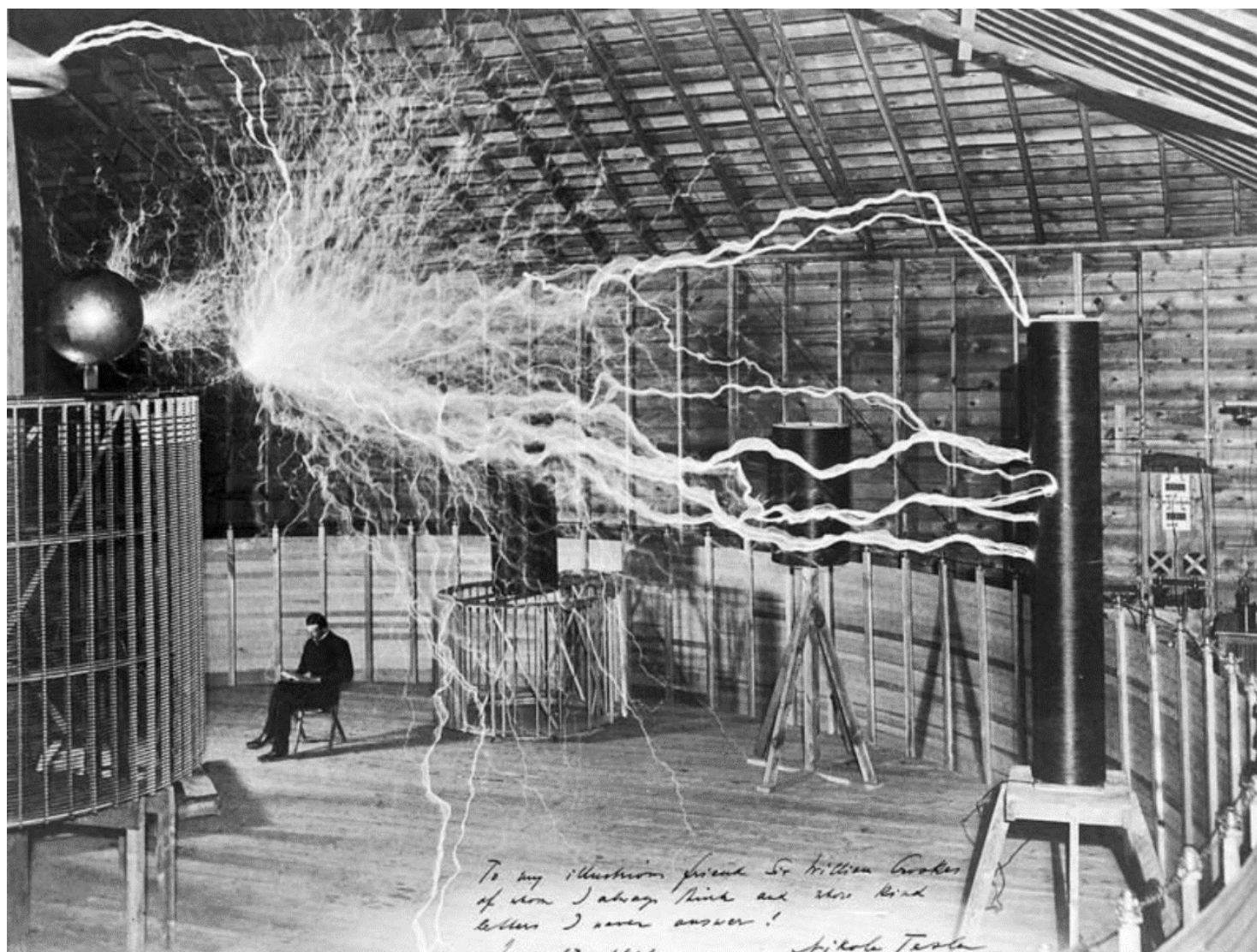
M	Topic & Assignment																								
M2	<div><h3>Titanic data mining analysis</h3><div><div>A. Background and overviews</div><div><ul style="list-style-type: none"><li><a href="https://www.rdocumentation.org/packages/titanic/versions/0.1.0">https://www.rdocumentation.org/packages/titanic/versions/0.1.0</a></li><li><a href="https://www.kaggle.com/competitions/titanic/overview">https://www.kaggle.com/competitions/titanic/overview</a></li><li><a href="https://www.encyclopedia-titanica.org/">https://www.encyclopedia-titanica.org/</a></li></ul><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 &amp; 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></div></div><div><div>B. Interesting models - built in R code for display convenience</div><div><pre>&gt; data &lt;- read.csv('titanic.csv')</pre><ul style="list-style-type: none"><li># Linear regression model</li><li>model &lt;- lm(survived ~ age + sex + pclass + sibsp + parch, data = data)</li></ul><div><li>Binomial Predicting survival based on age, sex, and passenger class</li><li>model &lt;- glm(survived ~ age + sex + pclass, data = titanic, family = binomial)</li></div><div><li>Poisson - Predicting the count of siblings/spouses based on passenger age</li><li>model &lt;- glm(sibsp ~ age, data = titanic, family = poisson) summary(model)</li></div><div><li>Neg.Binomial - Predict count of parents/children by passenger age and sex</li><li>model &lt;- glm.nb(parch ~ age + sex, data = titanic) summary(model)</li></div></div></div><div><div>C. Data &lt;class.github&gt;</div><div><ul style="list-style-type: none"><li>raw data; unsplit and preprocessed [source: <a href="https://hbiostat.org/data/">https://hbiostat.org/data/</a> &lt;titanic.3&gt;</li><li>train, test; from kaggle</li></ul></div></div><div><div>D. Data dictionary</div><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>&lt;parent.child&gt; #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 &lt;Cherbourg, Southampton, and Queenstown&gt;</td></tr></table><div>References:<ol style="list-style-type: none"><li><a href="#">Harrell Jr, F.E.,(2002)</a>. Titanic data, Vanderbilt biostatistics <a href="#">datasets</a>. Vanderbilt University. Retrieved from: <a href="https://hbiostat.org/data/repo/titanic.html">https://hbiostat.org/data/repo/titanic.html</a>. Retrieved on 05.15.2023.</li></ol></div></div></div>	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>																								

Data preprocessing and supervised, unsupervised algorithm purpose

# Data Preprocessing



Unsupervised learning	Dimensionality Reduction	Supervised learning
Input data is unlabeled	<ul style="list-style-type: none"><li>• Feature Elicitation</li><li>• Meaningful Compression</li><li>• Structure Discovery</li><li>• Big data visualization</li></ul>	Input data is labeled
Has no feedback mechanism	Clustering <ul style="list-style-type: none"><li>• Recommender Systems</li><li>• Targeted Marketing</li><li>• Customer Segmentation</li></ul>	Has a feedback mechanism
Assigns properties of given data to classify it	Classification <ul style="list-style-type: none"><li>• Identity Fraud Detection</li><li>• Image Classification</li><li>• Customer Retention</li><li>• Diagnostics</li></ul>	Data is classified based on the training dataset
Divided into Clustering & Association	Regression <ul style="list-style-type: none"><li>• Population Growth Prediction</li><li>• Estimating life expectancy</li><li>• Market Forecasting</li><li>• Weather Forecasting</li><li>• Advertising Popularity Prediction</li></ul>	Divided into Regression & Classification
Used for analysis	<ul style="list-style-type: none"><li>• Real-time decisions</li><li>• Game AI</li><li>• Robot Navigation</li><li>• Learning Tasks</li><li>• Skill Acquisition</li></ul>	Used for prediction
Algorithms include: k-means clustering, hierarchical clustering, apriori algorithm		Algorithms include: decision trees, logistic regressions, support vector machine
A unknown number of classes		A known number of classes
		



To my illustrious friend Dr William Crookes  
of whom I always think and who had  
letters I never answer!  
June 17, 1891 Nikola Tesla

# Master templates

my.header.2

Mod	Topic & Assignment
2	

my.header.1

Wk	Focus & Medium	Weekly Topic & Assignment
0		

Wk	Weekly Topic & Assignment
<div>11</div> <div>Mar</div>	<div>Header:</div> <div>1) What is<div>1) Perhaps one of th</div></div> <div>2) What isn't<div>a) Getting your n</div></div> <div>3) The mechanics and process<ul style="list-style-type: none"><li>• Orient</li></ul></div> <div>4) problem<div>Templated techniques help you quickly</div></div>

Wk	Weekly Topic & Assignment
<div>11</div> <div>Mar</div>	<div>Templated writing techniques hel</div> <div>Use kernel sentences: simple, declarative, active sentences (N.Chomsky)<div>Use of clear and concise language that is free of jargon and technical terms focuses the reader.<div>a) Joh</div></div></div> <div>1. Template: how.to. abc<div>1.1. item.1: Tai</div><div>1.2. item.s: U</div><div>1.3. item.: Us</div><div>1.4. item.: Ack</div><div>1.5. item.5: abdc</div></div> <div>Scenario: The</div>