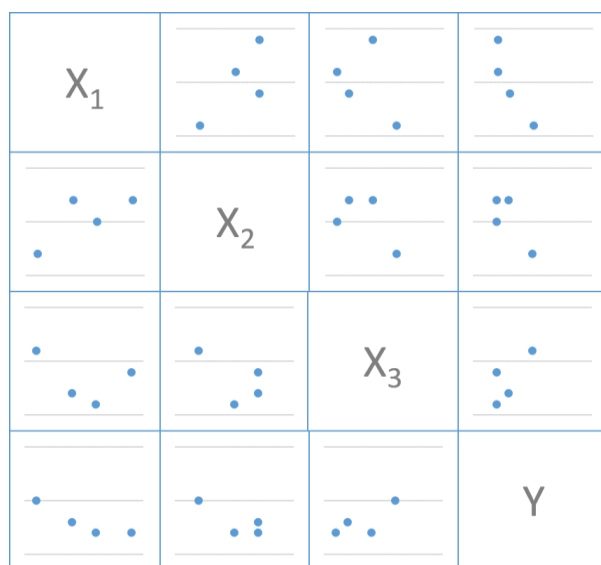


Data Science Essentials

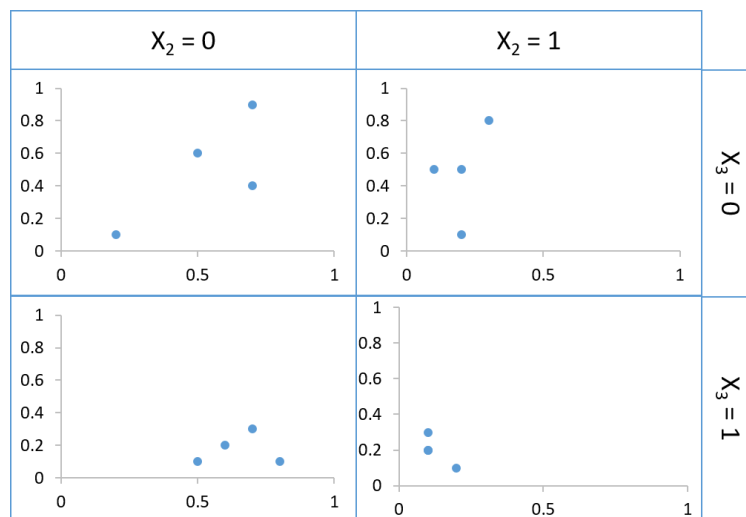
Visualizing Data

Data visualization is a highly useful way to explore data, and can help you determine apparent relationships between columns in order to identify candidates for predictive features in a machine learning model.

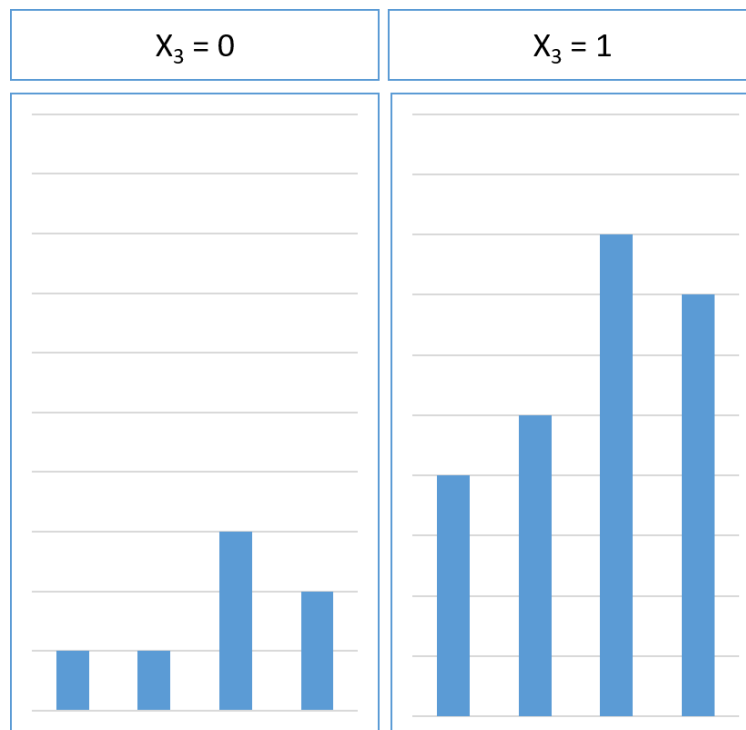
A scatter plot matrix, like the one below, shows scatter plots of selected columns in relation to each other, and is often a good starting point for data exploration. With a scatter plot matrix, you can easily spot variables that are collinear; which often indicates redundant features that should be removed from the model.



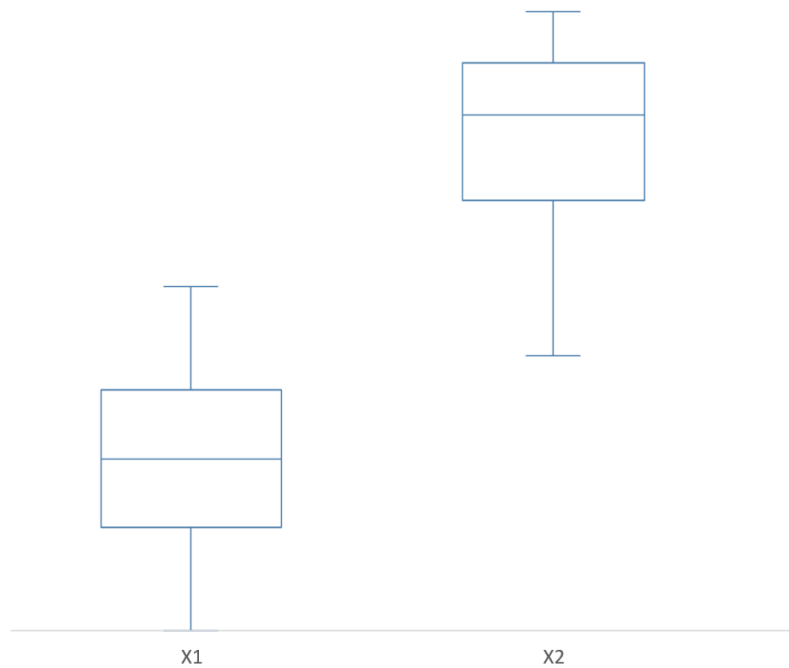
Scatter plots of individual columns can be useful for detailed exploration of the features in your dataset. A scatter plot enables you to see the intersection of values for two columns as plots in a chart. Additionally, you can condition the visualization on further columns; enabling you to visualize multiple dimensions of your data on a two dimensional chart. Scatter plots are particularly useful for spotting linear and non-linear relationships between variables.



Histograms that show the distribution of data density are useful for identifying potential outlier values. When conditioned, they show clearly how values for one variable may be distributed differently against specific values of another variable.



Box plots show the quartiles of numeric variables, with the median value indicated within a box that shows the first and third quartile. Additionally, lines known as *whiskers* can be extended from the box to show values outwith the values in the box. By comparing two values with box plots, you can easily see where the majority of the values for each variable lie, and to what extent the values in the two variables overlap.



You can generate visualizations using R or Python. See the following resources for more information.

- R Visualization Resources
 - Documentation for ggplot2: <http://docs.ggplot2.org/current/>
 - Cheat sheet for ggplot2: <https://www.rstudio.com/wp-content/uploads/2015/03/ggplot2-cheatsheet.pdf>
- Python Visualization Resources
 - Pandas plotting tutorial: <http://pandas.pydata.org/pandas-docs/stable/visualization.html>
 - Matplotlib tutorial: http://matplotlib.org/users/pyplot_tutorial.html