# What is Correlation?

Correlation is a measure of how strongly one variable depends on another.

- Correlation is **Positive** when the values **increase** together, and
- Correlation is **Negative** when one value **decreases** as the other increases
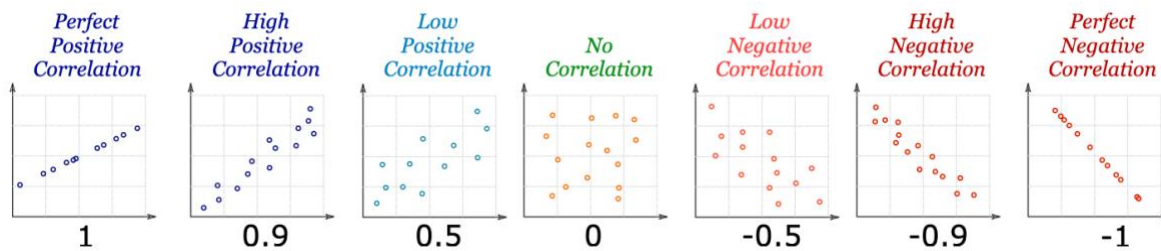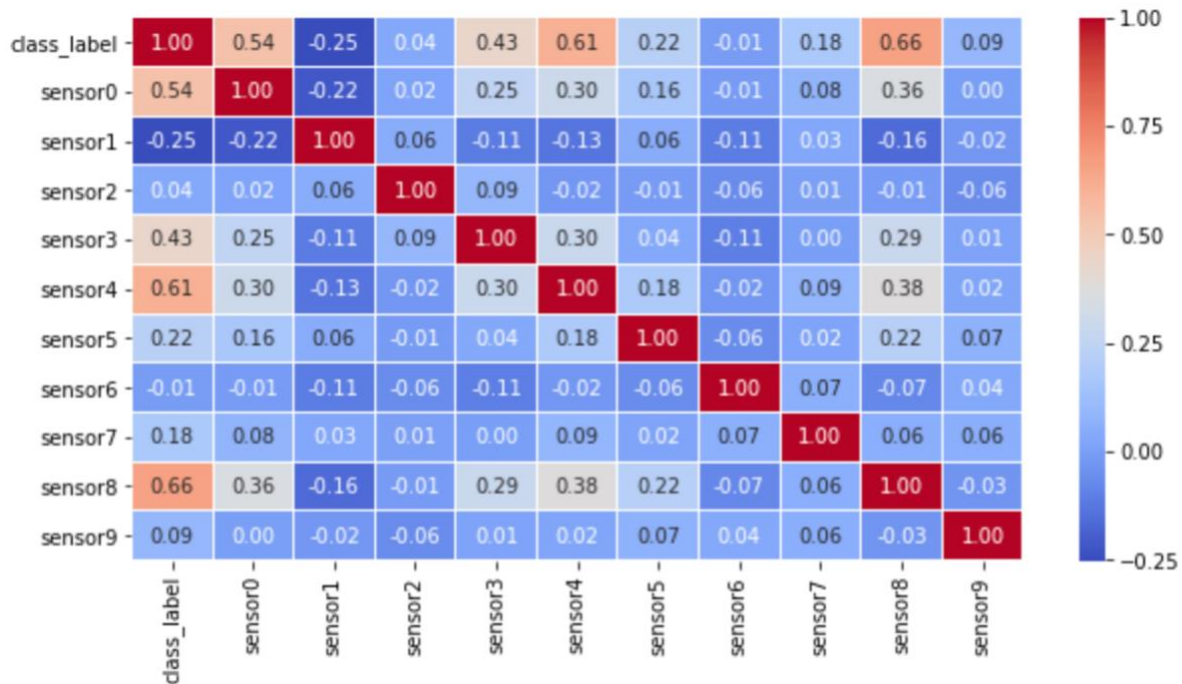- So, correlation varies between -1 and 1



Image Credit: https://www.mathsisfun.com/data/correlation.html

**Remember that Correlation is not causation**

"**Correlation is not causation**" means that just because two things **correlate does not** necessarily **mean** that one causes the other

The classic example of correlation not equaling causation can be found with ice cream and murder. That is, the rates of violent crime and murder have been known to jump when ice cream sales do. So, indeed buying ice cream and murders are highly correlated but it doesn't mean that one causes the other.

**Here is an example of Correlation Matrix:**

We can infer that class_label is strongly positively correlated with sensor 8 and has low correlation with Sensor 2 and 6



**Side-notes:**

**What is Multicollinearity?**

**Multicollinearity** refers to a situation in which two or more explanatory variables in a multiple regression model are highly linearly related. We have perfect **multicollinearity** if, for example, the correlation between two independent variables is equal to 1 or −1.

One of the assumptions of Classic Vanilla Multiple Regression is no multicollinearity. **Multiple regression** assumes that the independent variables are not highly correlated with each other.

**Why is Multicollinearity a problem with Multiple Regression?**

A simple way to look at it is to consider the following linear equation:

$$Y = mX + nZ + c$$

here m and n are coefficients of X and Z respectively; here coefficient m is increase in Y for every unit increase in X while holding Z constant. However, in practice it is impossible to hold Z constant. So, the positive correlation between X and Z mean that a unit increase in X is accompanied with some increase in Z at the same time.

So, it might be a good idea to remove highly correlated variables when you perform Linear Regression.

However, algorithms like Regression with Regularization or Decision Trees have feature selection embedded in them. So, it is not a bad idea to keep all the features and then let these kinds of non-parametric models or regularization techniques handle multicollinearity.

**What is the difference between parametric and non-parametric models?**

**Nonparametric statistics** is the branch of statistics that is not based solely on parametrized families of probability distributions (common examples of parameters are the mean and variance). Nonparametric statistics is based on either being distribution-free or having a specified distribution but with the distribution's parameters unspecified. The term parametric refers to parameters that define the distribution of the data.

Examples of parametric models are Logistic Regression, Naïve Bayes while Decision Trees doesn't make any assumptions regarding the distribution of the data hence, they are non-parametric.