Exploratory Data Analysis

- O Rectangular data is the busic data structure for statistical and machine learning mades
- OA coumn within a table is commonly referred to as feature.
- of now within a table is commonly referred to as a record.

Estimates of Location

Meon

Formula for the mean is pretty simple:

$$0 = \frac{1}{\sum_{i=1}^{n} x_i}$$

but if the values are sorted, and we want to trim the smallest and the largest values, we can use trimmed mean:

$$0 = \frac{\sum_{i=p+1}^{n-p} \times (i)}{n-2p}$$

Therefore climinating the influence of extreme values.

There is also weighted mean which is used when there are values that are much more variable only coming from a particular resource.

For example, if we are taking the overage from multiple sensors, and if one sensor is particularly variable, we may give the values collected from that sensor less weight.

Also, sometimes data does not represent the different groups we're interest—ed in equally. To correct that, we give a higher weight to the values from the group that over under represented.

$$0 \quad \overline{X}_{\omega} = \frac{\sum_{i=1}^{n} \omega_{i} \cdot x_{i}}{\sum_{i=1}^{n} \omega_{i}}$$

Outliers

Median is referred to as robust estimate of location since it is not influenced by outliers (extreme cases).

Outliers are not inherently invalid or erroneous. In contrast, outliers are sometimes informative in anomaly detection.

Still, outliers are often result of data errors such as mixing units (km vs m) or land readings from the source.

When outliers are result of a bad dato, mean will result in a poor estimate while the median will still be valid.

The median is not the only robust estimate of location. A trimmed mean con also be used to avoid the influence of outliers. It can be thought of as a compromise between the mean and the median.

Location in statistics is used to describe the central tendency of a dataset.

Estimoles of Voridoility

o Interquartile range is the difference between the 75th perantile and the 25th perantile. Also known IDR.

Mean absolute deviation

Let's say we have a dataset {1.4.41. The deviations from the mean are the differences: 1-3=-2 4-3=1 4-3=1

But we con't estimate the variability with only these values. In fact, the sum of the deviations from the mean is precisely zero.

To solve this, we take the absolute value of the deviations when we are using mean absolute deviation.

o mad = $\frac{\sum_{i=1}^{n} | \times_i - \overline{\times}|}{n}$ where $\overline{\times}$ is the sample mean.

The best known estimates of variability are the variance and the stan-

O Volicace =
$$S^2 = \frac{\sum_{i=1}^{n} (x_i - \overline{x})^2}{n-1}$$

· Standard deviation = 5 = IVaiance

You might worder, "Uny are we dividing by n-2". That's because, if you have a dotaget consist of 10 elements, sample of that dates t is naturally smaller. And excepting becomes much more closer to mean because of the small size.

It cases underestimation of the variance, so we use (n-1) instead of n to move the variance bigger.

Why (-1) because we only lose one degree of treedom when we estimate the variance of a sample, that is the mean of the sample.

O Median absolute deviation = median (|x,-ml, |x2-ml,..., |x,-ml)

A'= {12-51,14-71,15-51,16-51,110-51} = {0,1,1,3,5}

Median absolute deviation of A is the modular of A' -> 1

Estimates Based on Percentile

Statistics based on sorted data are referred to as order statistics.

The most basic measure is the range: the difference between the largest and smallest numbers.

In a detailet, Pth percentile means that at least (100-P) percent of the values take on this value or more.

Also, quantile is essentially the same as percentile, with quantiles indexed by fractions (so the .8 quantile is the same as the 80th percentile)

A common measurement of voridoility is the difference between 25th percentile and the 75th percentile, called the interquentile range or 16R

0 {3,1,5,7,6,7,2,9} → we sert these → {1,1,3,3,5,6,7,9} → 25th percentle is of 2.5 and the 75th percentle 6,5 so the intergration range is 6.5 - 2.5 = 4.

Exploring the Data Distribution

- D Location and varidoility are referred to as the first and second moments of a distribution.
- 0 The third and fourth moments are skewness and burtosis.
- D Stewness refer to whether the data is showed to larger or smaller values. (i.e if the shewness > 0, the data is showed to the right)
- O Kurtosis indicales the propensity of the data to have extrere values. (i.e high kurtosis means more data in the tails, low kurtosis means tever extreme values)
- · Skewness and leurtosis are gareally discovered through visualizations such as plots.

Exploring Binary and Categorical Data

- O When we have numeric values, there are often too many unique values to easily visualize, so we group numbers into loins.
- When we create bins, we actually create an ordered factor, which is note rally ordered. And we can create a histogram with this categorised data.

o like with bins, we can use something similar with categorical data which is bar chart.

Made, is the value, or values in case of a tie, that appears most often in the data.

Expected Value

- Of special type of categorical data is data in which the categories represent or can be mapped to discrete values on the some scale.
- · For example, let's think of a service that his three different subscription tiers.

For the sake of our scences, this service offers free webines. And after these webiners; 2.5 of the attendees will sign up for the \$300 service, 15% will sign up for the \$500 subscription and 2080 will not sign up for anything.

We can find the EV by basically estimating the weighted man of this detest

• EV is a fundamental concept in business valuation and copidal budgeting.

Correlation

$$0 = \frac{\sum_{i=1}^{n} (x_i - \overline{x})(y_i - \overline{y})}{(n-1) s_x s_y}$$
(Pearson's correlation coefficient)

The correlation coefficient always lies between +1 and -1, which respectively means perfect positive correlation and perfect negative correlation.

... More on the related notebode

Exploring Two or More Varidoles · Estimators like mean and variance can only look at one variable at a time, this is called univariate analysis. · Additionally, correlation analysis is an important method that compares two variables, this is called bivariate and yers. . In this section, we will look at additional estimates and plots, and at more than two voriables, that is called multivariate analysis. ... More on the related notebook