

Date.....

11/05/2020

PART 3: Practical Statistics

LESSON 1: DESCRIPTIVE STATISTICS PART - 1

1. Data Types

→ Two types of data:

a) Quantitative : data takes on numeric values that allows us to perform mathematical operations (like the number

of dogs).

b) Categorical : → These are used to label a

group or set of items (like dog

breeds like collies, labs, Poodles etc.).

→ This is also of two types:

a) Categorical Ordinal :

data take on a ranked ordering (like a ranked

interaction on a scale from

very poor to very good with

the dogs).

b) Categorical Nominal :

data do not have an order

or ranking (like the breeds

of the dog).

1.) Types of Quantitative data:

a) Continuous:

→ Data can be split into smaller & smaller units exists.

→ An example of this is the age of the dog - we can measure the units of the age in years, months, days, hours, seconds but there are smaller units that can be associated with the age.

b) Discrete:

→ Data only takes on countable values.

→ The number of dogs we interact with within a day is an example of a discrete data with two main types: Discrete & continuous.

2. Analyzing Quantitative Data

→ Four aspects for Quantitative data:

1. Measure of center.

2. Measure of spread.

3. The shape of the data.

4. Outliers.

3. Analyzing Categorical data

→ Analyzing categorical data has fewer parts to consider.

→ Categorical data is analyzed usually by looking at the counts or proportions of individuals that fall into each group.

Date.....

→ for example, if we were looking at the breeds of the dogs, we would care about how many dogs are of each breed, or what proportions of dogs are of each breed type.

a) Measures of center:

There are 3 measures of center:

1) Mean

2) Median

3) Mode

a) The Mean:

The mean is often called the average of the expected value in mathematics.

We calculate the mean by adding all of our values together, & dividing by the number of values in our dataset.

b) The Median:

- The median splits our data so that 50% of our values are lower & 50% are higher.

- for finding media, we must SORT the values in order.

Median for odd values:

If we have an odd no. of observ., the median is simply the middle number in the direct middle.

Date.....

- Median for Even values: If we have an even no. of observ., then the median is the average of the two middle values in the middle.

Ex: $\{1, 2, 3, 4, 5, 6, 7, 8\}$ → Median = $(5 + 6) / 2 = 5.5$

c) The Mode:

- The mode is the most frequently observed value in our dataset.

- There might be multiple modes for a particular dataset, or no mode at all.

• No Mode:

Ex: $\{1, 2, 3, 4, 5, 6, 7, 8\}$ → If all observ. in our dataset are observed with same freq., there is no mode.

Ex: $\{1, 2, 3, 4, 5, 6, 7, 8, 9\}$

• Many modes:

→ If two (or more) numbers share the maximum value, then there

is more than one mode.

Ex: $\{1, 2, 3, 3, 4, 5, 6, 6, 7, 8, 9\}$.

Here, 2 modes $\rightarrow 3 \& 6$.

2. Notation:

→ Notation is a common language used to communicate mathematical ideas.

→ "Think of notation as a universal language used by academic & industry professionals to convey mathematical ideas".

Date.....

→ Notation has the following properties:

To understand how to correctly use notation makes you seem really smart.

o) It allows you to read documentation, and implement our idea to your own problem.

(b) It makes ideas that are hard to say in plain words easier to convey.

iii) Random Variable:

→ Notated by a capital letter X .

• Observed value: shown as x_1, x_2, \dots
Notated by lowercase letters
($x_1 = 5$)

LESSON 2: DESCRIPTIVE STATISTICS PART - 2

Q) Measures of Spread:

→ measures of spread are used to provide us an idea of how spread out data are from one another.

→ Common measures of spread include:

1) Range

2) Interquartile Range (IQR)

3) Standard deviation

4) Variance

Date.....

2. Histograms:

→ Histograms are superuseful for understanding the different aspects of quantitative data.

3. Calculating the 5 number summary:

→ The 5 number summary consist of 5 values:

1. Minimum: the smallest number in the dataset.

2. Q₁: The value such that 25% of the data fall below.

3. Q₂: The value such that 50% of the data fall below.

4. Q₃: The value such that 75% of the data fall below.

5. Maximum: the largest value in the dataset.

→ Range:

The range is then calculated as the difference b/w the maximum and the minimum.

→ IQR:

The interquartile range is calculated as the difference b/w Q₃ & Q₁.

⇒ Box plot:

→ Boxplot is useful for quickly comparing & highlighting the spread of two datasets.

⇒ Standard Deviation and Variance:

→ The standard deviation is one of the most common measures for talking about the spread of data.

→ It is defined as "the average distance of each observation from the mean".

$$\text{Standard deviation} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$$

→ Variance:

Finding the average squared difference

with respect of each observation from the mean.

$$\text{Variance} = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2$$

⇒ Important final points

The variance is used to compare the spread of two different groups. A set of data with higher variance is more spread out than a dataset with lower variance.

Be careful though there might be an outlier (or outliers) that is increasing the variance, when most of the data are actually very close.

Date.....

2. When comparing the spread of two datasets, the units of each must be the same and must have the same scale and orientation.
 3. When data is related to money or the economy, higher variance (or standard deviation) is associated with higher risk.
 4. The standard deviation is used more often than the variance, because it shares the units of the original dataset.
3. Shape: \rightarrow from a histogram we can quickly identify the shape of our data, which helps influence all of the measures we learned in previous concepts.
- \rightarrow We learned that the distribution of our data is usually frequently associated with one of the 8 shapes:
- 1) Right-skewed \rightarrow mean > median
 - 2) Left-skewed \rightarrow mean < median
 - 3) Symmetric (frequently normally distributed)
 \rightarrow mean = median.

4) Outliers: ~~they will have a large influence on the mean and standard deviation.~~

→ Outliers have points that fall very far from the rest of our data points.

→ This influences measures like the mean or standard deviation much more than measures associated with the five number summary.

• common Techniques

When outliers are present, we should consider the following points:

1. Noticing they exist & their impact on summary statistics.
2. If type - remove or fix.
3. Understanding why they exist and what their impact on questions we are trying to answer about our data.

4. Reporting the 5 number summary which even though it is often a better indication of the data than measures like the mean & standard deviation when we have outliers.

5. Be careful in reporting. know how to ask the right questions.

• Descriptive Statistics

- Descriptive statistics is about describing our collected data.

• Inferential Statistics

- Inferential statistics is about using our collected data to draw conclusions to a larger population.