

Summary measures

Preliminaries

- Σ is shorthand for addition
- Suppose x_i is the *ith* observation.

$$x_1 + x_2 + x_3 = \sum_{i=1}^3 x_i$$

- If a is a constant

$$\sum_{i=1}^n a = na$$

- Mixed example:

$$\sum_{i=1}^n ax_i = a \sum_{i=1}^n x_i$$

Summary measures

Preliminaries

- b a constant
- n observations
- Simplify $\sum_{i=1}^n (a + bx_i)$
- Group terms:

$$\begin{aligned}\sum_{i=1}^n (a + bx_i) &= a + bx_1 + a + bx_2 + \dots + a + bx_n \\ &= a + a + a \dots bx_1 + bx_2 + \dots bx_n \\ &= an + b \sum_{i=1}^n x_i\end{aligned}$$

Summary measures

Preliminaries

- Let $\{x_i, y_i\}$ be paired observations
- k another constant.

- Simplify:

$$\sum_{i=1}^n (ak + bx_i + c(x_i y_i)) = ???$$

- Treat $x_i y_i$ as any other variable that varies in i .

$$\sum_{i=1}^n (ak + bx_i + c(x_i y_i)) = nak + \sum_{i=1}^n bx_i + \sum_{i=1}^n (cx_i y_i)$$

- Pull out the constants:

$$nak + \sum_{i=1}^n bx_i + \sum_{i=1}^n (cx_i y_i) = nak + b \sum_{i=1}^n x_i + c \sum_{i=1}^n (x_i y_i)$$

Summary measures

Statistics that describe distributions

- Measures of central tendency

- ① Mean

- ② Median

- ③ Mode

- Measures of variation

- ① Range

- ② Variance

- ③ Standard Deviation

- Measures of shape

- ① "Skew"

Summary measures

Measures of central tendency

- **Mean**

- Most common
- Central tendency

$$\hat{\mu}_x = \frac{1}{n} \sum_{i=1}^n x_i$$

- **Median**

- Mid-point of the data.
- To calculate:

- ① Order the data.
- ② Calculate $(n + 1)/2$
- ③ *Take the value at $(n + 1)/2$, or the average of the two closest (if $(n + 1)/2$ is not whole)*

Summary measures

Measures of central tendency (cont.)

- **Mode**

- The value that occurs most often

- Question: which measure(s) of central tendency are not affected by extreme values?

⇒ *Median and Mode*

Summary measures

Measures of variation

- **Range**

- A measure of data dispersion, though not used for many applications
- To calculate:
 - ① Identify the largest observation
 - ② Identify the smallest observation
 - ③ Take the difference

Summary measures

Measures of variation (cont.)

- **Sample Variance**

- The most commonly used measure of dispersion.
- Summarizes the how far a typical observation is from the mean

$$\hat{\sigma}_x^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_x)^2$$

- Why do we divide by $n - 1$ instead of n ?

⇒ We used one "piece" of information to calculate $\hat{\mu}_x$

Summary measures

Measures of variation (cont.)

- **Sample Standard deviation**

$$\hat{\sigma}_x = \sqrt{\hat{\sigma}_x^2}$$

- This is more desirable than the sample variance. Why?
 - Same scale as the mean.

Covariance

Relationships

- Covariance describes the relationship between two random variables

$$\hat{\sigma}_{xy} = \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_x)(y_i - \hat{\mu}_y)$$

- When will covariance be positive/negative?
 - $\hat{\sigma}_{xy} > 0 \Rightarrow$ tends to have $x_i > \hat{\mu}_x$ when $y_i > \hat{\mu}_y$
(and vice versa)
 - $\hat{\sigma}_{xy} < 0 \Rightarrow$ tends to have $x_i > \hat{\mu}_x$ when $y_i < \hat{\mu}_y$
(and vice versa)
- Covariance describes a "linear" relationship
- Any non-linear relationships with zero covariance?

Correlation

Basic

- Correlation describes a linear relationship
- $\hat{\rho}_{xy} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x \hat{\sigma}_y}$
- $\hat{\rho}_{xy} \in [-1, 1]$
- Can you prove that correlation is between -1 and 1?

Data Scaling

Mean

- What happens when we scale variables by either adding a constant or multiplying by a constant?
- Sample: $X = \{2, 3, 4\}$
- Calculate the mean of X :

$$\begin{aligned}\hat{\mu}_x &= \frac{1}{n} \sum_{i=1}^n x_i \\ &= \frac{1}{3} (2 + 3 + 4) \\ &= 3\end{aligned}$$

- Calculate the variance of X :

$$\begin{aligned}\hat{\sigma}_x^2 &= \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_x)^2 \\ \hat{\sigma}_x^2 &= \frac{1}{2} \left((2-3)^2 + (3-3)^2 + (4-3)^2 \right) = \frac{1}{2} (1 + 0 + 1) = 1\end{aligned}$$

Scaling

Adding a constant

- What if we define a new variable, $Z = X + 3$
- Calculate the mean of $Z = \{5, 6, 7\}$
- What will happen to the mean, variance?

$$\begin{aligned}\hat{\mu}_z &= \frac{1}{3} (5 + 6 + 7) = 6 \\ \hat{\sigma}_z^2 &= \frac{1}{2} \left((5 - 6)^2 + (6 - 6)^2 + (7 - 6)^2 \right) \\ &= \frac{1}{2} (1 + 0 + 1) = 1\end{aligned}$$

- Adding a constant **does** affect central tendency.
- Adding a constant **does not** affect dispersion.

- 1 **2** **3** **4** 5 6 7 8 9
- 1 2 3 4 **5** **6** **7** 8 9

Scaling

Multiplying by a constant

- What if we define a new variable, $J = 3X$
- Calculate the mean/variance of $J = \{6, 9, 12\}$
- What will happen to the mean, variance?

$$\hat{\mu}_z = \frac{1}{3} (6 + 9 + 12) = 9$$

$$\begin{aligned}\hat{\sigma}_z^2 &= \hat{\sigma}_j^2 = \frac{1}{2} \left((6 - 9)^2 + (9 - 9)^2 + (12 - 9)^2 \right) \\ &= \frac{1}{2} (9 + 0 + 9) = 9\end{aligned}$$

- Multiplying by a constant affects both mean and variance.

1 **2** **3** **4** 5 6 7 8 9 10 11 12

1 2 3 4 5 **6** 7 8 **9** 10 11 **12**

- Generally, if $J = aX$, then $\hat{\sigma}_j^2 = a^2 \hat{\sigma}_x^2$, $\hat{\mu}_j = a \hat{\mu}_x$

Covariance

Scaling

- Suppose $Z = aX$. What is $\hat{\sigma}_{zy}$?
- Write $\hat{\sigma}_{zy}$

$$\hat{\sigma}_{zy} = \frac{1}{n-1} \sum_{i=1}^n (z_i - \hat{\mu}_z) (y_i - \hat{\mu}_y)$$

- Substitute for $\hat{\mu}_z$ and z_i

$$\hat{\sigma}_{zy} = \frac{1}{n-1} \sum_{i=1}^n (ax_i - a\hat{\mu}_x) (y_i - \hat{\mu}_y)$$

- Factor out a and simplify:

$$\hat{\sigma}_{zy} = a \frac{1}{n-1} \sum_{i=1}^n (x_i - \hat{\mu}_x) (y_i - \hat{\mu}_y) = a\hat{\sigma}_{xy}$$

- Covariance is sensitive to scale!! Is this a problem?

Correlation

Scaling

- Is correlation sensitive to scale?

- Suppose $Z = aX$, $a > 0$.

- $$\hat{\rho}_{zy} = \frac{\hat{\sigma}_{zy}}{\hat{\sigma}_z \hat{\sigma}_y} = \frac{a \hat{\sigma}_{xy}}{a \hat{\sigma}_x \hat{\sigma}_y} = \frac{\hat{\sigma}_{xy}}{\hat{\sigma}_x \hat{\sigma}_y} = \hat{\rho}_{xy}$$