# Statistics 2 : Sampling and Covariance

Anand Systla

Masters in Financial Engineering Bootcamp
UCLA Anderson

August 16, 2022

# Today's Agenda

- Discuss the previous class take home questions ✗

- Clarifications

- Moments of Normal Distribution

- Sample vs Population

- Covariance and Correlation

- Correlation and Causation

# CLARIFICATIONS

- Poisson intensity $\lambda$ vs Exponential intensity $\beta = \left(\frac{1}{\lambda}\right)$ Units become important!

  (1) ▶ If $\lambda = 3$ (say 3 buses an hour), if we use $X \sim Exp(\beta = \frac{1}{\lambda} = \frac{1}{3})$

  $$P(0 \leq X \leq 1) = \int_0^{1\ hour} f(x) = \int_0^1 \lambda e^{-\lambda x} dx = \left(-e^{-\lambda x}\right]_0^1 = 1 - e^{-3} \approx 0.95 = F(1) - F(0)$$

  (2) ▶ If $\lambda = 3/60 = 1/20$ (say 1 bus every 20 min), if we use $X \sim Exp(\beta = \frac{1}{\lambda} = 20)$

  $$P(0 \leq X \leq 1) = \int_0^{1\ min} f(x) = \int_0^1 \lambda e^{-\lambda x} dx = \left[-e^{-\lambda x}\right]_0^1 = 1 - e^{-1/20} \approx 0.05 = F(1) - F(0)$$

  (3) ▶ Probability of seeing a bus in 30 minutes, $\lambda = 3$

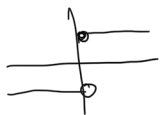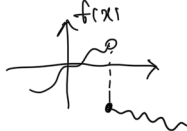  $$Y \sim Pois(\lambda_{60} = 3 \iff \lambda_{30} = 3/2) \implies P(Y = 1) = \frac{e^{-3/2}(3/2)^1}{1!} = 0.33$$

  $$X \sim Exp(\beta = \left(\frac{1}{\lambda_{60}}\right) = \frac{1}{3}) \implies P(X \leq \frac{1}{2}) = \left(\int_0^{1/2} f(x) dx = F(1/2) - F(0) = 1 - e^{-3 \cdot \frac{1}{2}} = 0.33\right)$$

- Modulus function $|x|$ properties
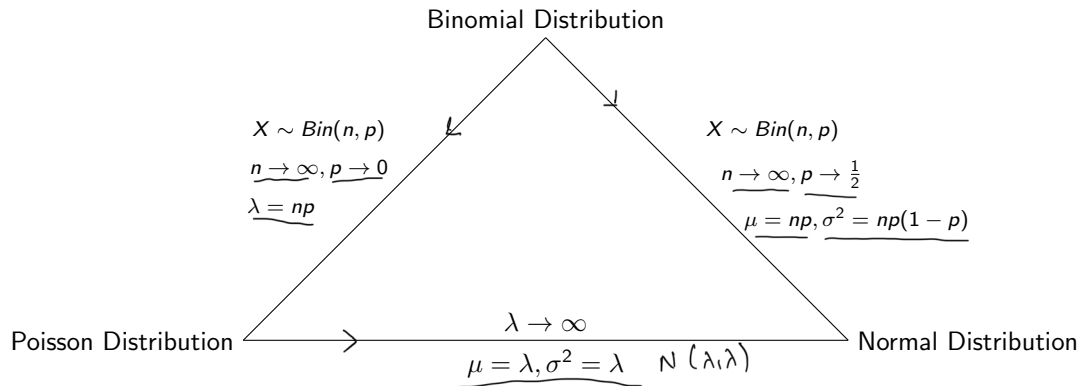
  ▶ $|x|$ is continuous everywhere but $|x|$ is not differentiable at 0. Easy to see this for the $x^2$ quadratic function $\lim_{x \to 0^+} 2x = \lim_{x \to 0^-} 2x = 0$

  $$\frac{d}{dx}|x| = \begin{cases} -1, & x < 0 \\ +1, & x \geq 0 \end{cases}, \qquad \lim_{x \to 0^+} \frac{d}{dx}|x| = 1 \quad \text{and} \quad \lim_{x \to 0^-} \frac{d}{dx}|x| = -1$$

# RELATIONSHIP BETWEEN DISTRIBUTIONS



Binomial Distribution
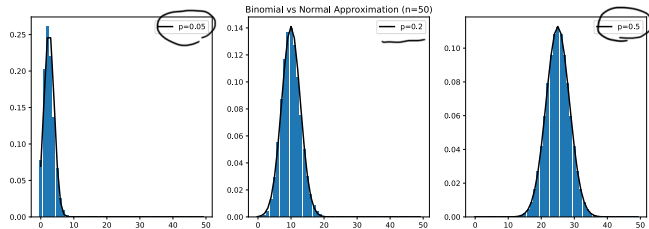
$X \sim Bin(n, p)$

$n \to \infty, p \to 0$

$\lambda = np$

$X \sim Bin(n, p)$

$n \to \infty, p \to \frac{1}{2}$

$\mu = np, \sigma^2 = np(1 - p)$

Poisson Distribution

$\lambda \to \infty$

$\mu = \lambda, \sigma^2 = \lambda$ $N(\lambda, \lambda)$

Normal Distribution

- A rough guideline to ensure the Normal approximation of the Binomial is reasonable
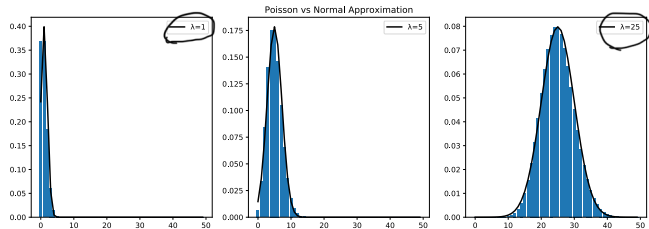  - ▶ $np \geq 10$ }
  - ▶ $n(1 - p) \geq 10$ }

$np \geq 10$ , $p \sim 0.1$, $n \geq 100$

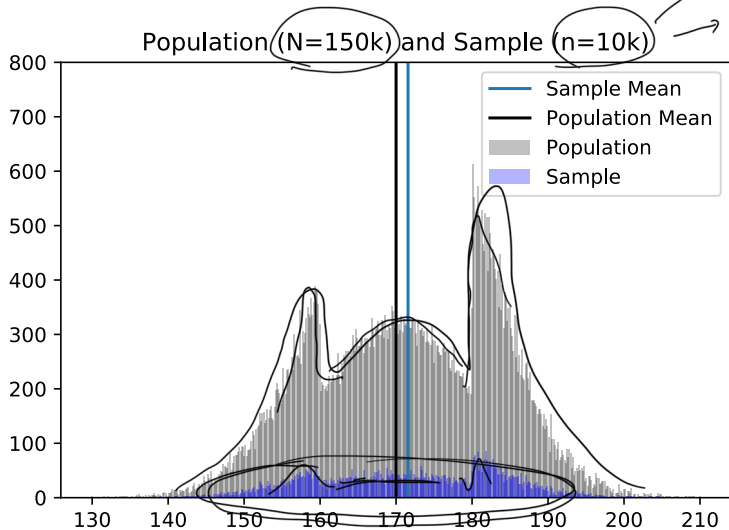# RELATIONSHIP BETWEEN DISTRIBUTIONS : EXAMPLES



(A)

# SAMPLE VS POPULATION

- Population includes all the data of a specified group. Sample is a subset of the population

|  | Population | Sampling Methodology |
|---|---|---|
| Height of people in US | 330 Mn | Selecting people from each state |
| Height of people in UCLA | 50,000 | Asking MFE/MBA/Professors |
| Weight of people in Japan | 125 Mn | Setting up volunteer "booths" in Tokyo |

- Gold standard is having a random sample that is representative of the population. Generally samples suffer from sampling/selection bias. In our case - (1) non-responsiveness, (2) under-coverage, (3) location of advertising

- Population is summarized by parameters. A sample is summarized by sample statistics. As the sample size approaches the population size, the sample statistic is going to approach population parameter    $(N)$

$$n \longrightarrow N$$

# POPULATION VS SAMPLE DATA



Population (N=150k) and Sample (n=10k)

Legend:
- Sample Mean
- Population Mean
- Population
- Sample

9k : $m$

9k : $m_2$

$m_k$

# Population Parameters and Sample Statistics

- Moments are robust ways of summarizing a RV. Common moments of interest are - mean, variance, skewness, kurtosis, quantiles, etc.

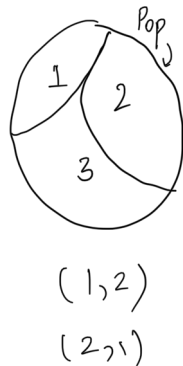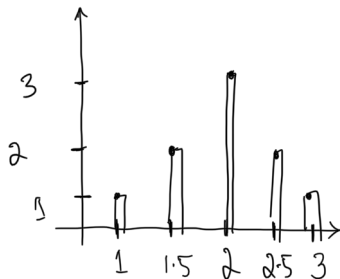| Moment | Population Parameter | Sample Statistic |
|--------|---------------------|------------------|
| Mean | $\mu = E[X] = \dfrac{\sum x_i}{N} = \int x f(x) dx$ | $\bar{X} = m = \dfrac{\sum x_i}{n}$ |
| Variance | $\sigma^2 = E[X^2] - E[X]^2 = \sum_i^N \dfrac{(x_i - \mu)^2}{N}$ | $s^2 = \sum_i^n \dfrac{(x_i - m)^2}{n-1}$ |

"Unbiased estimator"

# Population Parameter vs Sample Statistic : Example

- Sample data is different from sample statistic. Both can have their own distributions. An example is given below where population mean $\mu = \frac{1+2+3}{3} = 2$. Generating a sample of 2 draws with replacement and ordering matters

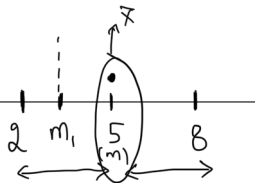| Samples | Mean ($\bar{X}$) |
|---------|------------------|
| (1,1)   | 1                |
| (1,2)   | 1.5              |
| (1,3)   | 2                |
| (2,1)   | 1.5              |
| (2,2)   | 2                |
| (2,3)   | 2.5              |
| (3,1)   | 2                |
| (3,2)   | 2.5              |
| (3,3)   | 3                |

Pop

$(1,2)$

$(2,1)$

- Does the distribution look familiar? More on this next class!

# SAMPLE VARIANCE

- We have seen that the sample variance is given by $s^2 = \sum_i^n \frac{(x_i - \bar{X})^2}{n-1}$

$(5-2)^2 + (5-8)^2 = 3^2 + 3^2 = 18$

$(4-2)^2 + (4-8)^2 = 4 + 16 = 20$



$s^2(\bar{X})$ , $\begin{array}{l} m = 5 \\ m_1 = 4 \end{array}$

- **Intuition 1 :** We use one data point in computing the mean. If we compute the sample mean $\bar{X}$ from $n$-data points, we no longer have $n$ independent data points. We can back out the $n$th number from $n-1$ data points and the mean.

- **Intuition 2:** The variance computes the squared deviation around the mean. $s^2$ computes variance centered around $\bar{X}$ and $\sigma^2$ computes the variance around $\mu$, which is different from $\bar{X}$. Any deviation from the sample mean $\bar{X}$ will only increase the varaince. So we bump up the sample variance by dividing by a smaller number $n-1$. Dividing by $n-1$ increases $s^2(\bar{X})$, closer towards $s^2(\mu)$.

- Say you are given $X = \{20\}$, what is $\bar{X} = 20$ and what is $s^2(\bar{X}) = NaN$

$$S^2(\bar{x}) = \frac{\sum (x_i - \bar{x})^2}{n-1} \leq \sum (x_i - \mu)^2$$

$n$ $\{1, 2, 3\}$ $\rightarrow$ $\{1, 2, 3, \bar{x}\}$

$$\bar{x} = \frac{6}{3} = 2$$

$$\underset{x_1}{1}, \underset{x_2}{2}, \frac{2}{\bar{x}}$$

$$\bar{x} = \frac{x_1 + x_2 + x_3}{3}$$

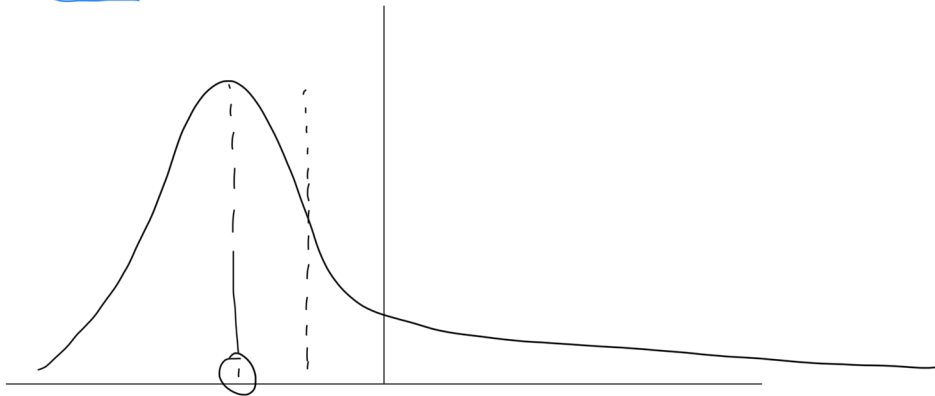$$\Rightarrow x_3 = 3\bar{x} - x_1 - x_2 = 6 - 1 - 2 = 3$$

# SKEWNESS AND KURTOSIS

$> 0$
$< 0$

- Skewness: Which direction is the data (tail) drawn out towards? $\frac{1}{N}\frac{\sum(x-\mu)^3}{\sigma^3}$

- Kurtosis: How much weight do the tails have? $\frac{1}{N}\frac{\sum(x-\mu)^4}{\sigma^4}$
  - "Peakedness" has nothing to do with the kurtosis. Two distributions can have the same mean and standard deviation, but can have different weights placed on their tails
  - Kurtosis $> 3$ is leptokurtic ($< 3$ is platykurtic)
  - Kurt($N(0,1)$) = 3, so people are often interested in excess kurtosis ($=$ kurtosis$-3$)

# COVARIANCE

- Lets us study joint variation of two random variables and how they co-move
- In our sample below $\bar{X} = 100$ and $\bar{Y} = 10$. The covariance asks are $X - \bar{X}$ and $Y - \bar{Y}$ above and below zero together?

| $X$ | $Y$ | $X - \bar{X}$ | $Y - \bar{Y}$ | $(X - \bar{X})(Y - \bar{Y})$ |
|-----|-----|---------------|---------------|------------------------------|
| 100 | 10  | 0             | 0             | 0                            |
| 102 | 9   | 2             | -1            | -2                           |
| 98  | 11  | -2            | 1             | -2                           |
| 110 | 14  | 10            | 4             | 40                           |
| 90  | 6   | -10           | -4            | 40                           |

$$E\left[(x-\bar{x})(y-\bar{y})\right]$$
$$= \frac{80}{5}$$

- Properties of covariance ~~terms~~
  - $cov(X, Y) = E[(X - \bar{X})(Y - \bar{Y})] = E[(X - \bar{X})Y] - E[(X - \bar{X})\bar{Y}] = E[Y(X - \bar{X})]$
  - $cov(X, Y) = E[XY - \bar{X}Y - X\bar{Y} - \bar{X}\bar{Y}] = E[XY] - E[X]E[Y]$
  - $cov(X, c) = 0$
  - $cov(aX, X) = a \cdot cov(X, X) = a \cdot var(X)$
  - $cov(X, Y + c) = cov(X, Y)$
  - $cov(X, Y + Z) = cov(X, Y) + cov(X, Z)$

# COVARIANCE VISUALIZATION

- $cov(X, Y) \lessgtr 0,$   $cov(X, X) > 0,$   outliers?

$X = \{100, 102, 98, \; 110, 90\}$

$Y = \{10, 9, 11, \; 14, 6\}$



- Covariance does not tell us how far the points are from the dotted line. It also does not tell us how steep or flat our line is

COV > 0

COV > 0

y

y

$\overline{x}$

# CORRELATION

$Z \sim N(0,1):$ $\quad X = \mu + \sigma Z \Rightarrow Z = \dfrac{X-\mu}{\sigma}$

- Although the covariance gives us a measure of co-movement of two random variables, it scales with any constant multiplying the random variable $cov(aX, Y) = a\,cov(X, Y)$

- Correlation does not depend on the scale of the data. It is a measure of linear dependence and $\rho_{XY} \in [-1, 1]$. Why?

$$\rho_{XY} = \underbrace{\frac{cov(X, Y)}{\sigma_X \sigma_Y}} = \underbrace{\frac{\overbrace{E[(X - \bar{X})(Y - \bar{Y})]}^{Unit1 \times Unit2}}{\underbrace{\sigma_X \sigma_Y}_{Unit1 \;\; Unit2}}} = E\left( \underbrace{\frac{X - \bar{X}}{\sigma_X}}_{\tilde{X}} \underbrace{\frac{Y - \bar{Y}}{\sigma_Y}}_{\tilde{Y}} \right)$$

- Properties of the correlation

  - ▶ Correlation is a dimensionless quantity
  - ▶ $corr(aX, Y) = corr(X, Y)$
  - ▶ $corr(X, Y + c) = corr(X, Y)$

    $cov(X, y) = E[XY] - E[X]E[Y]$

  - ▶ If $X$ and $Y$ are independent $\rho_{XY} = 0$, but $\rho_{XY} = 0$ does not imply independence (non-linearity)
  - ▶ If $X \sim N(0, 1)$, then $cov(X, X^2) = E[(x - \bar{x})(x^2 - \bar{x^2})] = E[X \cdot X^2] - E[X] \cdot E[X^2]$

    " dependence "

$$= E[x^3] - E[\cancel{X}] \cdot E[x^2] = 0$$

$$\begin{cases} corr(X, Y) = 0.75 \\ corr(Y, Z) = 0.25 \end{cases}$$

$$corr(X, Z) = ?(x)$$

$$\begin{array}{c c} & \begin{array}{c c c} X & Y & Z \end{array} \\ \begin{array}{c} X \\ Y \\ Z \end{array} & \left[ \begin{array}{c c c} 1 & P_{xy} & x \\ & 1 & P_{yz} \\ P_{xy} & & \\ x & P_{yz} & 1 \end{array} \right] \end{array}$$

$$det(corr) \geqslant 0$$

$$1(1 - P_{yz}^2) - P_{xy}(P_{xy} - x \cdot P_{yz}) + x(P_{xy} \cdot P_{yz} - x) \geqslant 0$$

$$a x^2 + bx + c \geqslant 0 \qquad x \in [P^L, P^H]$$

# CORRELATION VISUALIZATION

- Correlation can be 1 irrespective of how spread out of narrow the data is
- How does our confidence on the correlation measure change with number of data points?