

## STATISTICS CLASS 3

- Agenda for today's class
  - Correlation and causation
  - MGF
  - Central limit theorem
  - Law of large no:
  - ANOVA (One-way)

- Correlation and causation

$$\rho_{xy} = \text{corr}(x, y) = \frac{\text{cov}(x, y)}{\sigma_x \cdot \sigma_y}$$

$$\rho_{xy} = \frac{E[(x - \bar{x})(y - \bar{y})]}{\sigma_x \sigma_y} = E \left[ \overbrace{\left( \frac{x - \bar{x}}{\sigma_x} \right)}^{\tilde{x}} \cdot \overbrace{\left( \frac{y - \bar{y}}{\sigma_y} \right)}^{\tilde{y}} \right]$$

standardize first  
and then get covar.

[Cauchy] / r = 1/2

$$\left[ \begin{array}{l} \text{Schwarz} \\ \text{Inequality} \end{array} \right] (E[XY]) \leq E[X^2] \cdot E[Y^2]$$

Properties of correlation:

(a)  $\rho \in [-1, 1]$  dimensionless

(b) If  $X$  &  $Y$  are independent  
 $\text{cov}(X, Y) = 0, \rho_{XY} = 0$

But  $\text{cov}(X, Y) = \rho_{XY} = 0 \nRightarrow$  Independence

(c) Correl. only captures linear relationship

$$X \sim N(0, 1), Y = X^2 =$$

$$\text{cov}(X, Y) = \text{cov}(X, X^2) = E[X \cdot X^2] - E[X] \cdot E[X^2]$$

$$E[X^3] - E[X] E[X^2] = 0$$

$$(d) V[X+Y] = V[X] + V[Y] + 2 \cdot \text{cov}(X, Y)$$

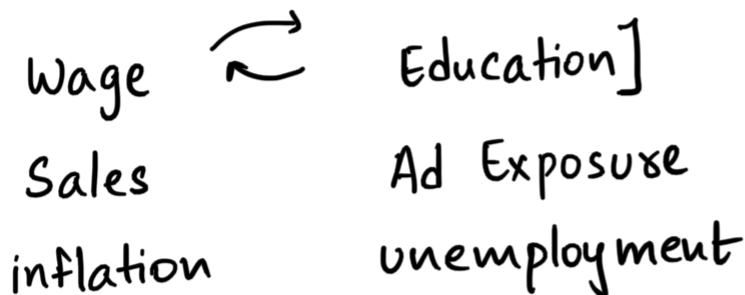
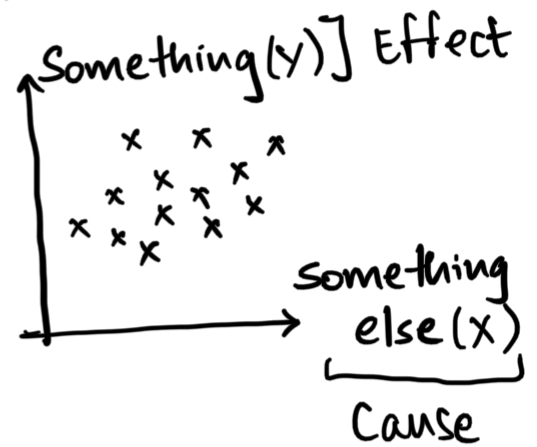
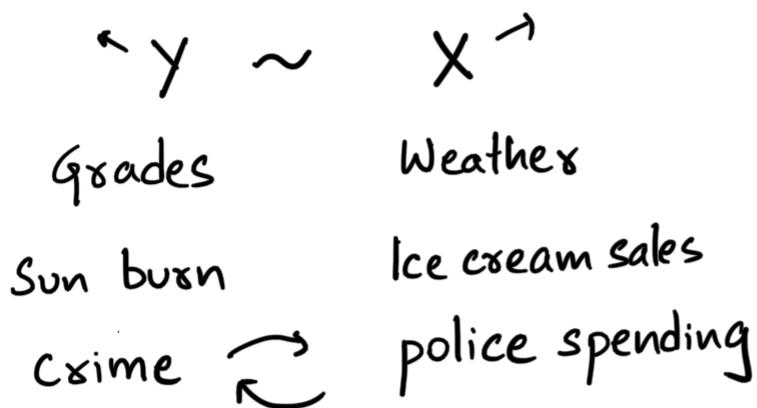
$$= \sigma_X^2 + \sigma_Y^2 + 2 \cdot \sigma_X \sigma_Y \cdot \rho_{XY}$$

$$= \text{corr}(X, Y) = \text{corr}(X, Y)$$

$$(e) \text{corr}(X, aY) = \text{corr}(X, Y)$$

$$(f) \text{corr}(X, Y+c) = \text{corr}(X, Y)$$

Studying causal relationships is the fundamental problem of econometric inference.



- Causation goes a step further than correl.  
change  $x \rightarrow$  change  $y$ ?

- Establishing causality
  - Experiment (gives control)
  - Cause occurs before effect
  - If X occurs, Y also occurs
  - No alternate explanations

- Moment Generating Function (MGF)

$$\mu = E[X] = \int x \cdot f_X(x) dx$$

$$\sigma^2 = V[X] = \underbrace{E[X^2]} - (E[X])^2$$

MGF for a RV X is defined as

$$\underbrace{M_X(t)}_{\downarrow} = E[e^{tX}] = \begin{cases} \sum e^{tk} \cdot P(X=k) \\ \int_{-\infty}^{\infty} e^{tx} \cdot f_X(x) \cdot dx \end{cases}$$

Why does the MGF work?

$$e^{tX} = 1 + tX + \frac{(tX)^2}{2!} + \frac{t^3}{3!} \cdot X^3 + \dots$$

$$\int_{-\infty}^{\infty} e^{tX} = 1 + t E[X] + \frac{1}{2} t^2 E[X^2] + \dots \quad \left( \frac{t^n}{n!} \right) = 1$$

$$[E[e^t]] = 1 + \frac{E[X]}{1!} t + \frac{E[X^2]}{2!} t^2 + \dots$$

$$\mu = \left. \frac{\partial}{\partial t} E[e^{tx}] \right|_{t=0} = E[X] + \frac{t}{1!} E[X^2] + \dots \Big|_{t=0} = E[X]$$

$$\left. \frac{\partial^2}{\partial t^2} E[t^2 x] \right|_{t=0} = E[X^2] + t \cdot E[X^3] + \dots \Big|_{t=0} = E[X^2]$$

Ex:  $X \sim \text{Binom}(n, p)$

$$M_X(t) = E[e^{tx}] = \sum_{x=0}^n e^{tx} \cdot \binom{n}{x} p^x (1-p)^{n-x} \Big|_{(x+y)^n = \sum_{k=0}^n \binom{n}{k} x^k y^{n-k}}$$

$$= \sum \binom{n}{x} (p \cdot e^t)^x (1-p)^{n-x} = (pe^t + q)^n$$

$$\mu = \left. \frac{\partial}{\partial t} M_X(t) \right|_{t=0} = \left. n \cdot (pe^t + q)^{n-1} \cdot pe^t \right|_{t=0} = n \cdot p$$

Properties of MGF:

(a) If  $M_X(t) = M_Y(t)$  then  $F_X(x) = F_Y(x)$

(b) If  $Y = X_1 + X_2 + \dots + X_n$  independent RVs

$$M_Y(t) = E[e^{ty}] = E[e^{t(X_1 + X_2 + \dots + X_n)}]$$

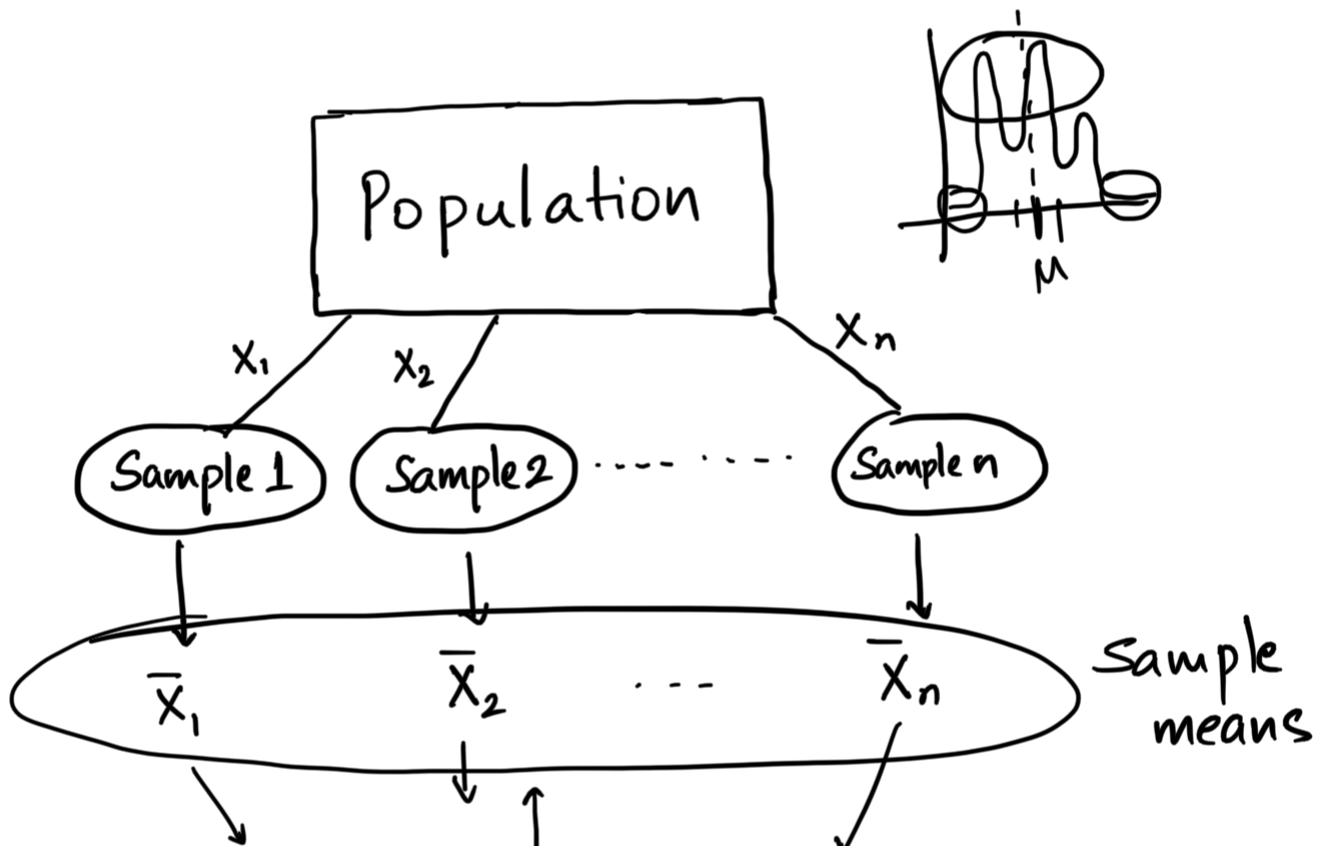
$$M_Y(t) = M_{X_1}(t) \cdot M_{X_2}(t) \cdots M_{X_n}(t)$$

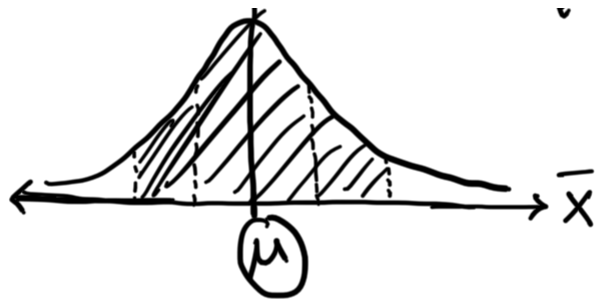
- Law of Large No: and Central Limit Thm.

If  $X_1, X_2, \dots$  are iid, from some population  $P$  with a defined mean and variance.

[LOLN]:  $\bar{X}_n \rightarrow \mu$  as  $n \rightarrow \infty$  with prob. 1

[CLT]:  $\sqrt{n} \left( \frac{\bar{X}_n - \mu}{\sigma} \right) \xrightarrow{d} N(0,1)$





95% confident  $\bar{X}_n \in \left[ \mu - 2\frac{\sigma}{\sqrt{n}}, \mu + 2\frac{\sigma}{\sqrt{n}} \right]$

(OR)  $\bar{X}_n \in \left[ \mu - 2\frac{\sigma}{\sqrt{n}}, \mu + 2\frac{\sigma}{\sqrt{n}} \right]$  with 95% probability

$$\underbrace{\bar{X}(n)} = \frac{\bar{X}_1 + \bar{X}_2 + \dots + \bar{X}_n}{n}$$

LOLN says  $\underbrace{\bar{X}(n)} \xrightarrow{P} \underbrace{\mu}$  as  $\underbrace{n \rightarrow \infty}$  with prob. 1

CLT says  $\underbrace{\bar{X}(n)} \sim \underbrace{N\left(\mu, \frac{\sigma^2}{n}\right)}$

$\bar{X}(n)$  converges to  $\mu$  at rate  $1/\sqrt{n}$

$$\underbrace{E[\bar{X}(n)]} = \frac{\overbrace{\mu + \mu + \dots + \mu}^{n \text{ times}}}{\underbrace{E[\bar{X}_1] + E[\bar{X}_2] + \dots + E[\bar{X}_n]}_{n \cdot \mu}} = \frac{n \cdot \mu}{n} = \underbrace{\mu}$$

$$\underbrace{V[\bar{X}(n)]} = \frac{\overbrace{V[\bar{X}_1] + V[\bar{X}_2] + \dots + V[\bar{X}_n]}^{n \text{ times}}}{\underbrace{n^2}} = \frac{n \cdot \sigma^2}{n^2} = \underbrace{\frac{\sigma^2}{n}}$$

$$E[\bar{X}(n)] = \mu, \text{ SD}[\bar{X}(n)] = \sigma/\sqrt{n}$$

$$(\text{Standardization}): Z = \frac{\bar{X}(n) - \mu}{\sigma/\sqrt{n}}$$

Why does central limit theorem work?

$$- \bar{X}_i = (\mu) + \sqrt{\epsilon_i} \text{ Error}$$

When we avg  $\bar{X}_i$  the  $\epsilon_i$ 's cancel

- Low prob. values are not drawn freq.

Applications  $\begin{cases} \nearrow \text{Hypothesis testing} \\ \rightarrow \text{Elections} \\ \searrow \text{Bootstrapping} \end{cases}$

## • Analysis of Variance (ANOVA)

ANOVA  $\Leftrightarrow$  Dummy variable regression

Total variance is total variance can



Intuition of ANOVA is that variation

be decomposed into within & across group variation.

$$\chi^2_{n-1} \sim S^2 = \frac{\sum (x_i - \bar{x})^2}{(n-1)}, \quad SST = \sum (x_i - \bar{x})^2$$

$$F(X) = \left\{ \begin{array}{ccc} \text{10\%} & \text{15\%} & \text{20\%} \\ \underbrace{1, 1, 4}_{C1} & \underbrace{4, 5, 6}_{C2} & \underbrace{7, 7, 10}_{C3} \end{array} \right\}$$

$$F(Y) = \left\{ \underbrace{3, 4, 2}_{C2-C1}, \underbrace{6, 6, 4}_{C3-C1} \right\} \quad \mu = \frac{25}{6} = 4.16$$

$SST = 8.77$

$$C2-C1 = \{3, 4, 2\} \quad C3-C1 = \{6, 6, 4\}$$

$$\mu_1 = 3$$

$$\mu_2 = 16/3$$

$$SS21 = (-1)^2 + 1^2 = 2$$

$$SS31 =$$

$$\left. \begin{array}{l} \mu_1 = \mu_2 = \mu_m \\ \mu_1 - \mu_m = \mu_2 - \mu_m \end{array} \right\}$$

Q ANOVA asks - is there a diff. b/w the 3 classes?

$$SS_{\text{Total}} = \text{SS}_{\text{Within}} + SS_{\text{Across}}$$

$$\bar{\mu} = \frac{45}{9} = 5, \quad SST = (-4)^2 \times 2 + 1^2 \times 2 + 0 + 1^2 + 2^2 \times 2 + 5^2 = 68$$

$$C_1 = \{1, 1, 4\}$$

$$C_2 = \{4, 5, 6\}$$

$$C_3 = \{7, 7, 10\}$$

$$\mu_1 = 2$$

$$SS_1 = (-1)^2 + (-1)^2 + 2^2$$

$$\mu_2 = 5$$

$$SS_2 = (-1)^2 + 0 + 1^2$$

$$\mu_3 = 8$$

$$SS_3 = 1^2 + 1^2 + 2^2$$

$$\left[ \begin{array}{ccc} SS_1 = 6 & SS_2 = 2 & SS_3 = 6 \end{array} \right] SS_W$$

$$SS_{\text{Across}} = \sum n_i (\bar{\mu}_i - \bar{\mu})^2$$

$$= 3(-3)^2 + 3 \times 0 + 3 \times 3^2 = 54$$

$$SS_{\text{Total}} = 68$$

$$SS_{\text{Within}} = 6 + 2 + 6 = 14$$

$$SS_{\text{Across}} = 54$$

$\rightarrow \chi^2$

$$\begin{aligned} \text{F-stat} &= \frac{\text{Mean Sq. Across}}{\text{Mean sq. Within}} = \frac{\text{'SSA'}/(k-1)}{\text{SSW}/(n-k)} \end{aligned}$$

$\downarrow$   $\chi^2$

$$H_0: \underbrace{M_1 \neq M_2 = \dots = M_n}_{=0} \quad F=0$$

$$\begin{aligned} &\rightarrow M_1 = M_2 = M_M \\ &\rightarrow M_1 - M_M = M_2 - M_M \end{aligned}$$

Ratio of 2  $\chi^2$  distributions (iid) has F-dist.

$$\underbrace{U_1 \sim \chi_n^2}, \underbrace{U_m \sim \chi_m^2}$$

$$F = \frac{U_1/n}{U_2/m} \sim \text{F-dist. with } \underbrace{n \text{ and } m \text{ deg. of freedom}}$$

$$F = \frac{54/3-1}{14/9-3} \approx \underline{11.5} \quad \left| \quad F=1 \right.$$

Larger F value, more variation across groups

(a) Random sampling from pop<sup>n</sup>  $(\sigma^2)$

SS Across << SS Within

$$\Rightarrow \underline{0 \leq F \leq 1}$$

(b) Ex: here would be comparing diff.  
portfolio returns!

How diff. are means of diff. portfolios?