# Statistics Class 2

- Agenda for today
    - Sample v. Population
    - Moments
    - Covariance
    - Correlation

- Sample v. Population

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

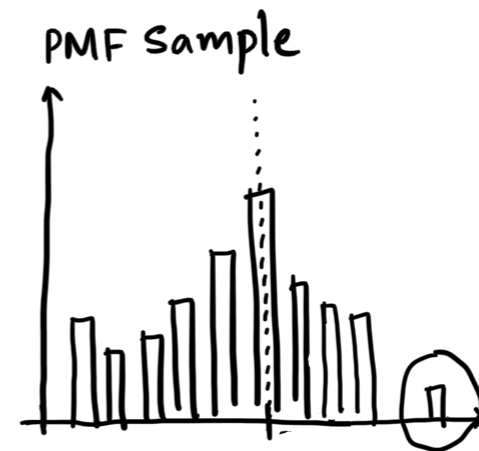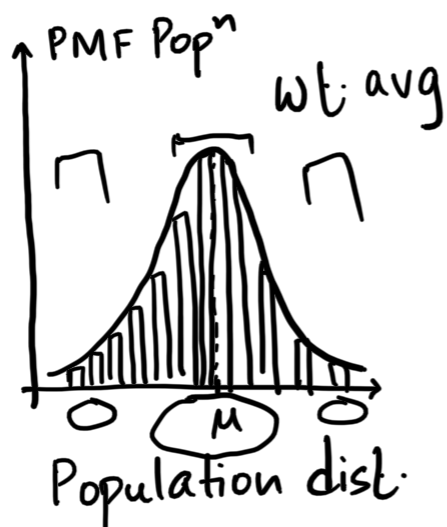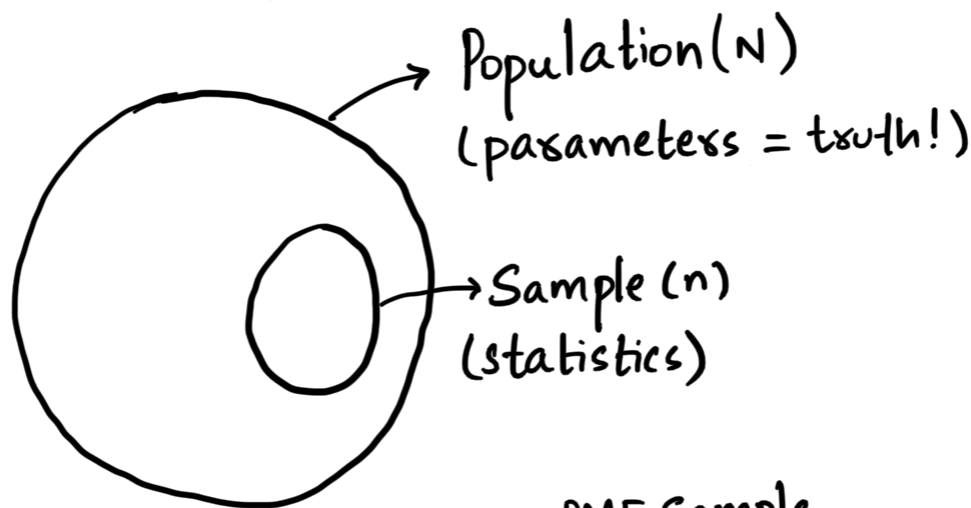  |  | Population | Sampling |
  |---|---|---|
  | Avg. Height of people in USA | ~330 Mn | select people from each state? |
  | Avg. Height of people in UCLA | ~50 k | MFE, MBA, Profs? |
  | Avg. Weight of people in Japan | ~125 Mn | Volunteering booths at diff |

places in Japan?

$$Sampling\ bias \begin{cases} \to Undercoverage\ (ex:2) \\ \to Advertising\ (ex:3) \\ \to Non\text{-}response\ (ex:1) \end{cases}$$

Random sample → gold standard
(good proxy/less bias)

Ex 1: Avg. Height of people in USA



→ Population (N)
(parameters = truth!)

→ Sample (n)
(statistics)

PMF Pop$^n$          wt; avg

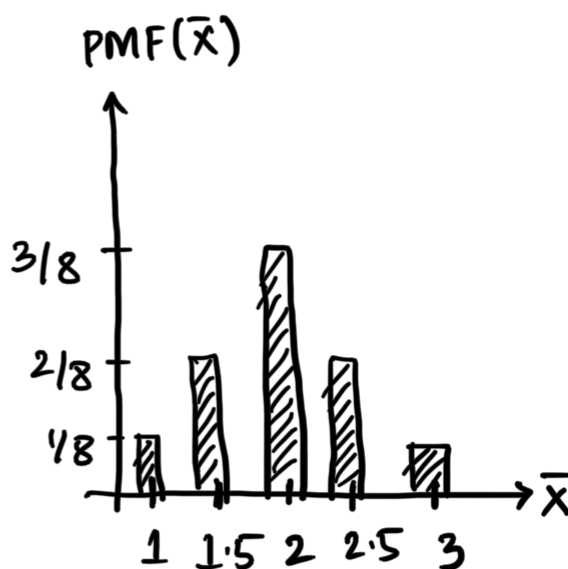Population dist.

PMF Sample

Distribution of sample data

# Ex 2: Distribution of sample statistic

Population mean $= \textcircled{M} = \dfrac{1+2+3}{3} = 2$

Sampling $= 2$ draws with replacement

Mean $(\bar{x}) \to$ proxies for $M$

| 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                |

PMF($\bar{x}$)



Depending on who samples from pop$^n$, $\bar{X}$ changes

Similarly we can get PDF of other sample stats.
- $\to$ Mean (Normal) ✓
- $\to$ Variance $(x^2)$
- $\to$ Median / Quantiles

$\searrow$ skewness/kurtosis

- Moments

  Robust ways of summarizing data

  $\{$ <u>Mean, var, std dev,</u> Quantile, <u>skew, kurt</u> $\}$

Sample data 1

$X = \{12, 14, 14, 20\}$



12   14   20

$\bar{X} = \dfrac{12 + 14 + 14 + 20}{4} = 15$

Sample data 2

$Y = \{15, 15, 15, 15\}$



15

$\bar{Y} = 15$

<u>Mean</u> $=$ Avg. distance from 0

$\searrow$ First moment

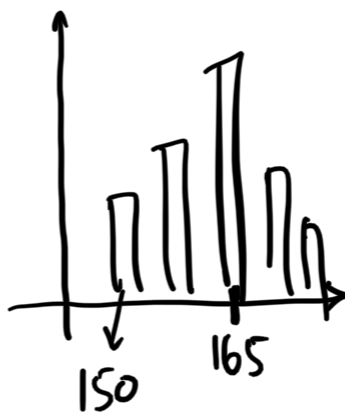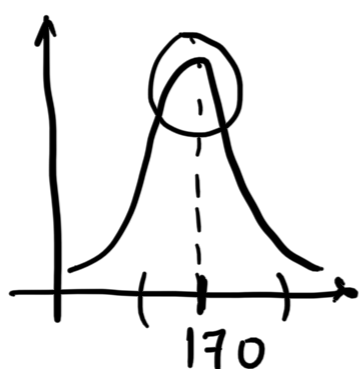$\left(\dfrac{1}{n} \sum x_i^2\right)$   $\dfrac{12^2 + 14^2 + 14^2 + 20^2}{4}$ $=$ 234

$\dfrac{1}{n} \sum y_i^2 = 15^2 = 225$

skewness ... more ...

$\rightsquigarrow$ squared dist. from $0$ (increases with more spread)

Alternatively we can compute dist. from mean

$$\underbrace{\frac{1}{n} \sum (x_i - \bar{x})^2}_{} = \frac{(-3)^2 + 2 \cdot (-1)^2 + (5)^2}{4} = 9 \qquad \frac{1}{n} \sum (y_i - \bar{y})^2 = 0$$

centered
$2^{nd}$ moment

$\downarrow$
No spread
in data

Remove effect of $1^{st}$ moment $(\bar{x})$ to capture any additional info. in the dataset.



170

150   165

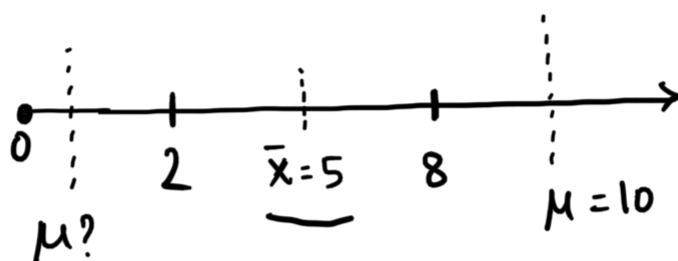| Moment | Population (param) | Sample (statistic) |
|---|---|---|
| (1) Mean | $\mu = E[x] = \int x_i \cdot \underbrace{f_x(x)}_{} \, dx$  $\mu = \sum x_i / N$ | $\bar{x} = m = \sum\limits_{i=1}^{n} x_i / n$ |

(2) Variance $\left[\sigma^2 = E[x^2] - E[x]\right.$ $\qquad$ $S^2 = \sum\limits_{i=1}^{} (x_i - \overset{m}{\underset{\downarrow}{\mu}})/n-1$

$\left. \sigma^2 = \sum\limits_{i=1}^{N} (x_i - \textcircled{\mu})^2/N \right.$

(3) skewness

(4) kurtosis

- Degrees of freedom

  - Intuition 1: Say $X = \{2, 8\}$



$0$    $2$   $\bar{X}=5$   $8$    $\mu = 10$

$\mu?$

$\dfrac{1}{n-1} \sum (x_i - \bar{X})^2 = (2-5)^2 + (8-5)^2 = 18$

$\dfrac{1}{n} \sum (x_i - \mu)^2 = (2-4)^2 + (8-4)^2 = 20 /\!/$

$\sum (x_i - 6)^2 = (2-6)^2 + (8-6)^2 = 20 /\!/$

Squared deviation at mean is smallest at $\bar{X}$

$\left[ S^2 \right.$ computes variance centered at $\bar{X}$

$\left[ \sigma^2 \right.$ computes variance centered at $\mu$

If $S^2(\bar{x}) = \frac{1}{n} \sum (x_i - \bar{x})^2$ } Function that takes $\bar{x}$ as input

then $\underline{S^2(\mu)} \geqslant \underline{S^2(\bar{x})}$
                                    lower bound

Division by $n-1$ ↑ $S^2(\bar{x})$ to be closer to $S^2(\mu)$

- Intuition 2:

Say $X \sim N(\mu, \sigma^2)$ } True data generation process

If $\mu = 0, \sigma = 1$

[Population]: $Z = \{ 0.5, 0.1, -0.2, 0.3, -1, 1.5 \}$

[Sample]:

| $X$ | $X - \bar{X}$ |
|---|---|
| 1.5 | 1.3 |
| 0.1 | -0.1 |
| -1 | -1.2 |

- If we know $\bar{X}$, one data point is redundant

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

free data points

$$\overline{x} = 0.2 \qquad 0 \qquad \Big| - \; \sqsupset n-1 \quad \dots$$

after computing the mean!

[DF] : No: of data points available to compute sample statistics
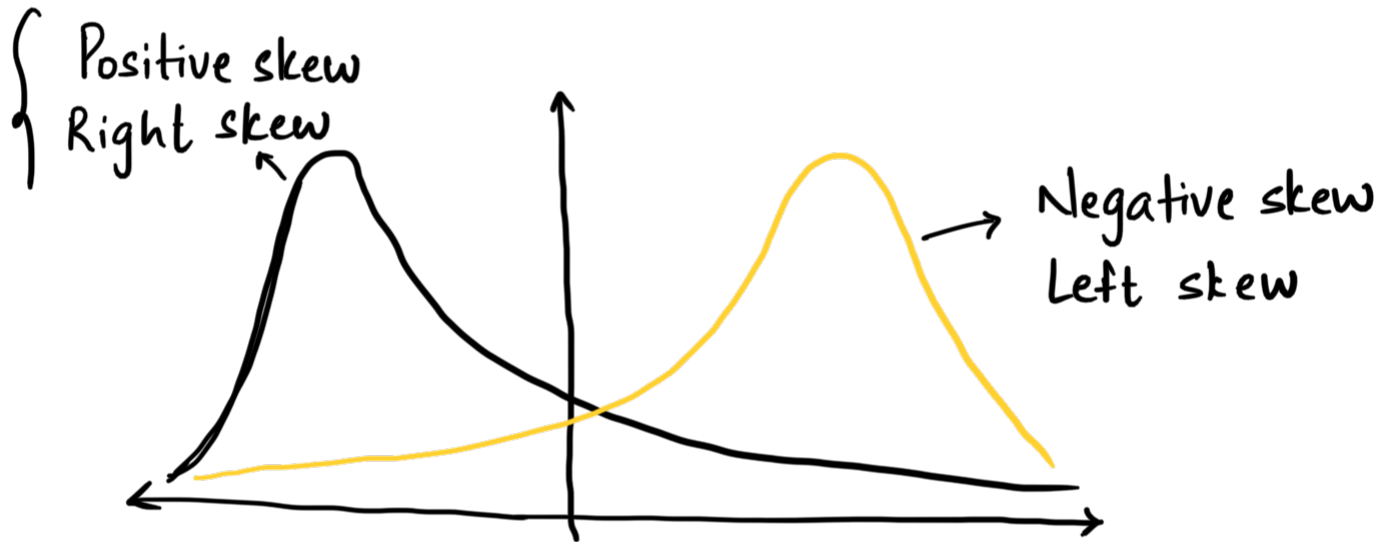
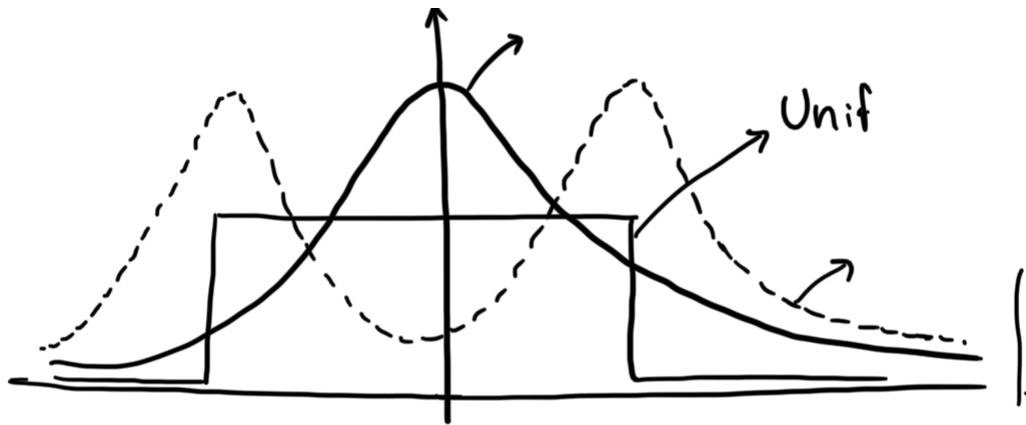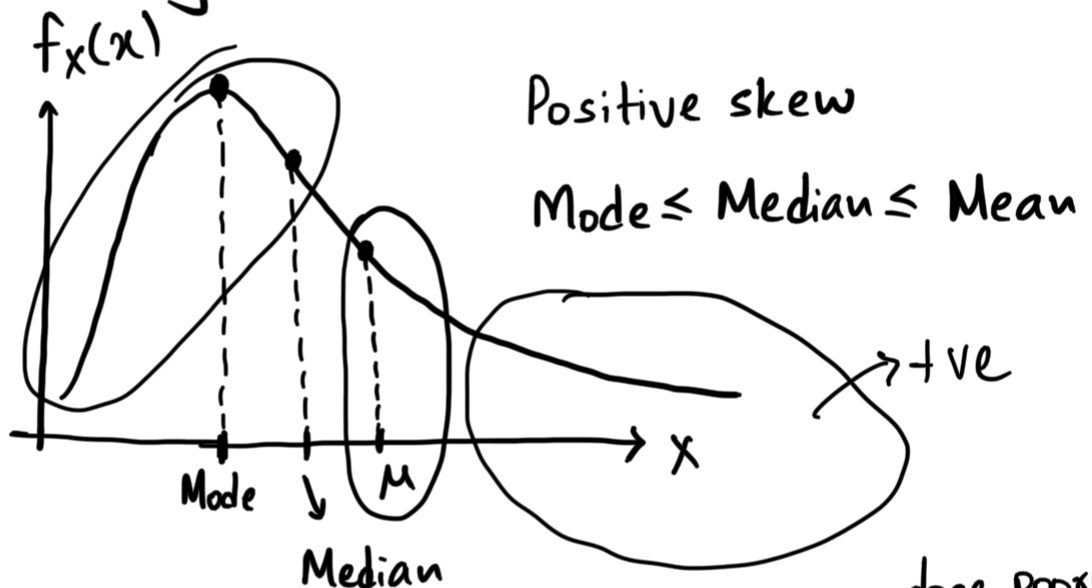| Sample stat | DF |
|---|---|
| (Mean) | $n$ |
| Var | $n-1$ |
| skew | $n-2$ |
| kurtosis | $n-3$ |
| $\vdots$ | $\vdots$ |

Ex:   $x = \{44\}$ sample

(a) $\overline{x} = 44$

(b) $s^2 = \dfrac{(44 - 44)}{\underbrace{n-1}_{=0}} = $ NaN

- Skewness

Unif

Positive skew
Right skew

Negative skew
Left skew

which way does the tail point?

$f_X(x)$

Positive skew

Mode ≤ Median ≤ Mean

+ve

Mode

M

Median

X

Mean - mode } does poorly on small

Pearsons mode skewness = $\dfrac{\rule{3cm}{0.4pt}}{\text{std. dev}}$ ⎰ samples

Pearson's median skewness = $\dfrac{(\text{Mean} - \text{Median}) \times 3}{\text{std. dev}}$

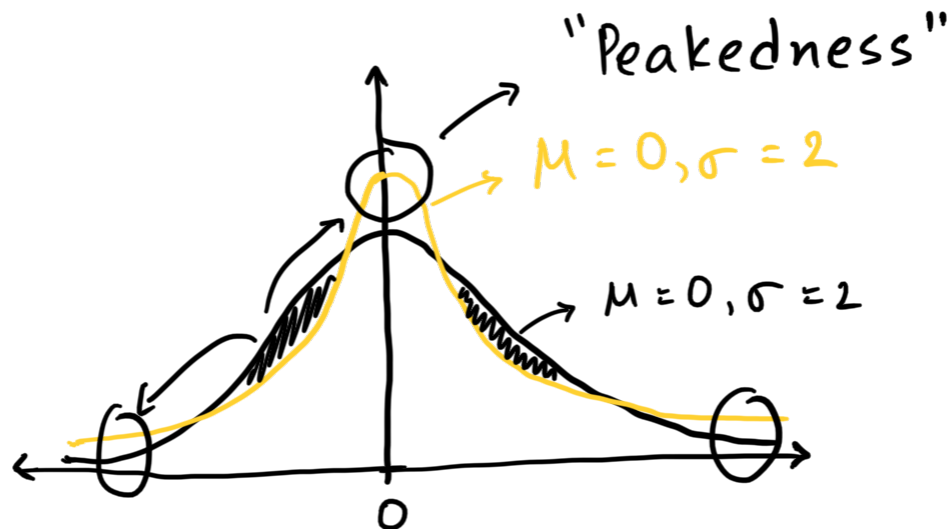Moment estim. skew = $\dfrac{1}{N} \underbrace{\dfrac{\sum (x_i - \mu)^3}{\sigma^3 /\!/}}_{\substack{\text{Removes effect of} \\ \text{mean and variance}}} = \dfrac{1}{N} \sum \left(\dfrac{x - \mu}{\sigma}\right)^{\!3}$

## $\underline{\text{Skew}\,(N(\cdot,\cdot)) = 0}$

- Kurtosis



"Peakedness"

$M=0, \sigma=2$

$M=0, \sigma=2$

Orange curve has fatter tails

kurtosis measures weight on tails

Moment estim. kurtosis $= \boxed{\left(\dfrac{1}{N} \ \dfrac{\Sigma(x_i - M)^4}{\sigma^4}\right) \mathbin{/\!/}}$

$\sqrt{\phantom{x}}$ kurtosis $(N(0,1)) = 3$, $\underline{\text{kurt} \ \varepsilon \ [1, \infty)}$

If kurtosis $> 3$ (leptokurtic)

kurtosis $< 3$ (platykurtic)

✓ Excess kurtosis $= \underbrace{kurt - 3}_{\text{sets Excess kurt. of std. normal} = 0}$

- Covariance

Captures joint variation of 2 RVs $X, Y$

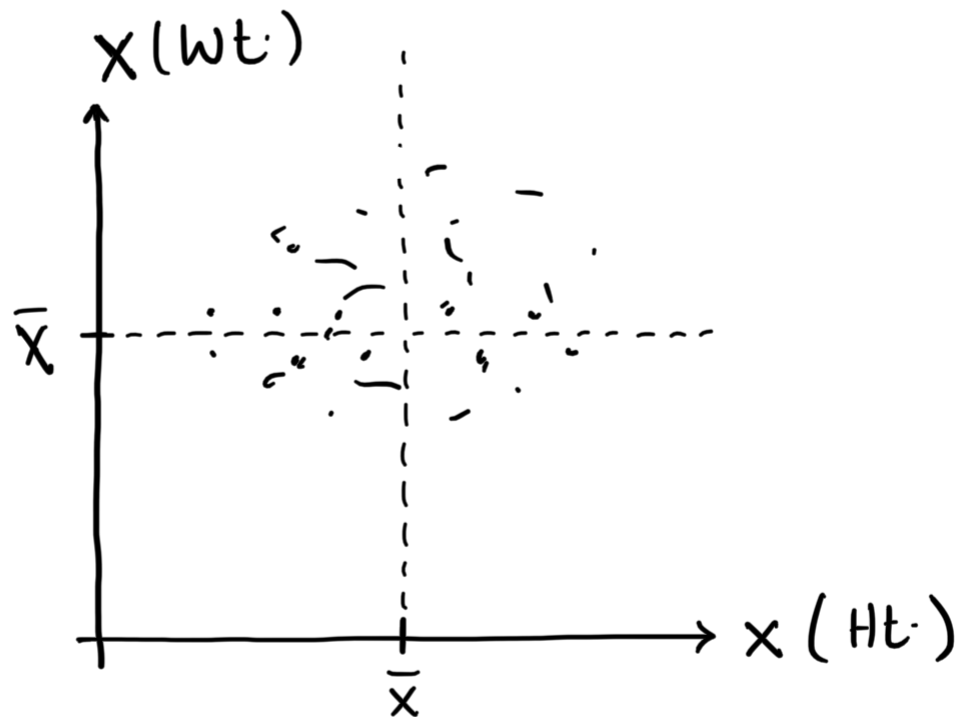| X | Y | $X - \bar{X}$ | $Y - \bar{Y}$ | $(X - \bar{X})(Y - \bar{Y})$ |
|---|---|---|---|---|
| 100 | 10 | 0 | 0 | 0 |
| 102 | 9 | 2 | $-1$ | $-2$ |
| 98 | 11 | $-2$ | 1 | $-2$ |
| 110 | 14 | 10 | 4 | 40 |
| 90 | 6 | $-10$ | $-4$ | 40 |

$\bar{X} = 100$ | $\bar{Y} = 10$

$$\text{Cov}(X, Y) = \frac{1}{-} \sum (x_i - \bar{X})(y_i - \bar{Y})$$

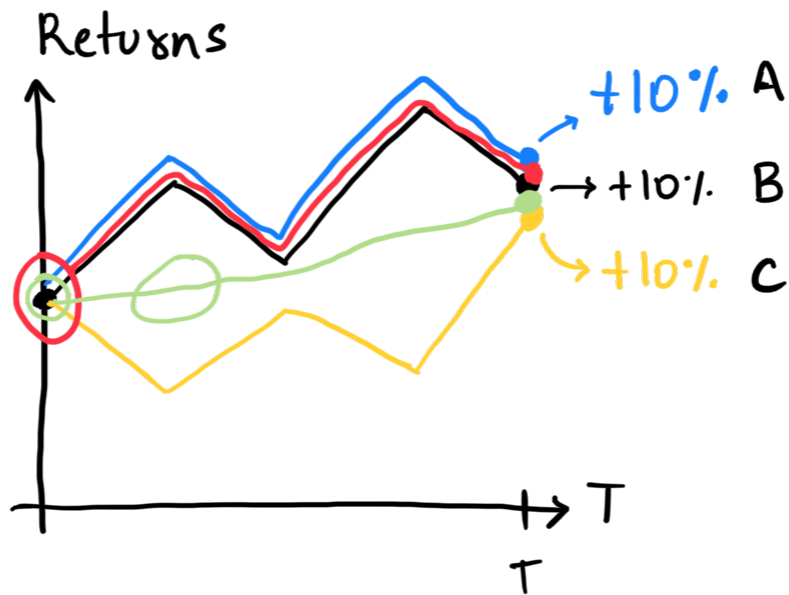$cov(X,Y) > 0$ : X,Y are above or below the mean together

$cov(X,Y) < 0$ : If X is above its mean, Y is below (vice-versa)

Ex 1:

X(wt.)

$\overline{X}$

$\underline{X}$

→ X (Ht.)

(a) Cov > 0

(b) Cov < 0

(c) cov(X,X)

(d) cov = 0

(e) Outlier

## Ex 2:

Returns

$+10\%$ A

$\to +10\%$ B

$\to +10\%$ C

T

$\cov(A,B) > 0$

$\cov(A,c) < 0$

$\cov(B,c) > 0$

"Markowitz" 1952

Portfolio 1: 50% A, 50% B

Portfolio 2: 25% A, 25% B, 50% C

---

Properties of covariance

(a) $\cov(x,y) = E[(x-\bar{x})(y-\bar{y})]\}$ const. 0

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

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

(b) $Cov(x,y) = E\left[xy - \bar{x}y - x\bar{y} + \bar{x}\bar{y}\right]$

$$\| = E[xy] - E[x]\cdot E[y]$$

$$\Rightarrow E[xy] = E[x]\cdot E[y] + Cov(x,y)$$

(c) $Cov(X,X) = Var(X)$

sensitivity to scale $\Big\{ Cov(X,aY) = a\cdot Cov(X,Y)$

$Cov(X,c) = 0$

$Cov(X,Y+c) = Cov(X,Y)$

$Cov(X+Y,Z) = Cov(X,Z) + Cov(Y,Z)$

$Cov(XY,Z) = E\left[(XY - E(XY))(Z - E[Z])\right]$

$$= E\left[z\left(XY - E[XY]\right)\right]$$

$$\neq cov(X,Z) \cdot cov(Y,Z)$$

(d) $Var(X+Y) = E\left[\left((X+Y - \bar{X} - \bar{Y})^2\right)\right]$

$$= E\left[\left(X+Y\right)^2 + (\bar{X}+\bar{Y})^2 - 2(X+Y)(\bar{X}+\bar{Y})\right]$$

$$= V[X] + V[Y] + 2 \cdot cov(X,Y)$$

- Correlation and Causation

$$\rho_{XY} = Corr(X,Y) = \frac{Cov(X,Y)}{\sigma_x \cdot \sigma_y} \left.\begin{array}{c}\} \text{ Scaled} \\ \text{covariance}\end{array}\right.$$

Properties of correlation

(a) $\rho_{xy} \in [-1, 1]$

scale
invar.
(b) $\text{Corr}(aX, Y) = \text{corr}(X, Y)$

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

Causality is the fundamental problem of econometric inference.

$$X \sim Y$$

| X | Y |
|---|---|
| Education | Wage |
| Ice cream sales | Sunburns |
| Drug | Disease |