

# Stats

- ① Covariance ✓
- ② Pearson Correlation Coefficient ✓ } → Homework ??
- ③ Spearman Rank Correlation Coefficient ✓ =
- ④ Log Normal Distribution ✓ = p value
- ⑤ Hypothesis Testing [ Null Hypothesis  
Alternate Hypothesis ] } Influential
- ⑥

## ① Covariance { Feature Selection }

X	Y	→ DATASET
→ 12	20	
→ 13	25	
→ 14	30	
→ 28	40	① Relationship between X & Y → features
→ 35	50	$\rightarrow \underline{X \uparrow} \quad y \uparrow \quad X \downarrow \quad y \uparrow$ $\quad \quad \quad X \uparrow \quad y \downarrow \quad \rightarrow \underline{X \downarrow} \quad y \downarrow$

② Numbers we can actually quantify the relationship →

Quantify → Covariance, Pearson Correlation,  
Spearman Rank Correlation

# ① Covariance

$n = \text{Sample}$

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

X	y
12	2
13	3

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

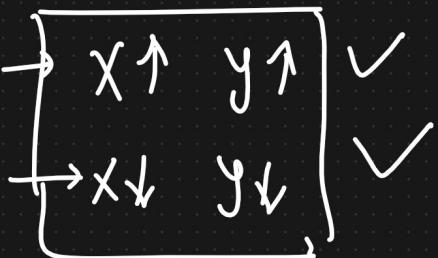
X	y
14	4

$\bar{x} = \text{mean of } x$

X	y
15	5

$$\bar{y} = \text{mean of } y = (12-13.5) * (2-3.5) + \\ (13-13.5) * (3-3.5) + (14-13.5) \\ * (4-3.5) + (15-13.5) * (5-3.5)$$

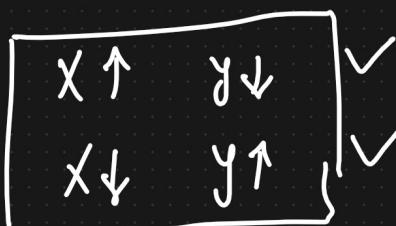
$n = \text{Sample}$



$$= (-1.5) * (-1.5) + (0.5) * (-0.5) + \\ = 1.667 // \rightarrow \text{vc} \quad > 100 \quad < 100 \\ 10,000 \quad 1000 \quad > 50 \quad < 50$$

X	y
5	2
4	3
3	4
2	5

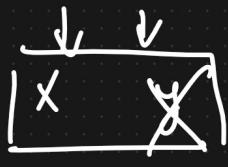
$$\text{Cov}(x, y) = \frac{-\text{vc}}{n-1}$$



$$\boxed{\text{Cov}(x, y) = 0} \quad ?.$$



Does this mean? ?



$$\text{Cov}(x, y) = 0$$

limiting the Quantified

result

Compare  
↔

1	1
1	1
1	1

① +ve or -ve { No limit }

{ limit }

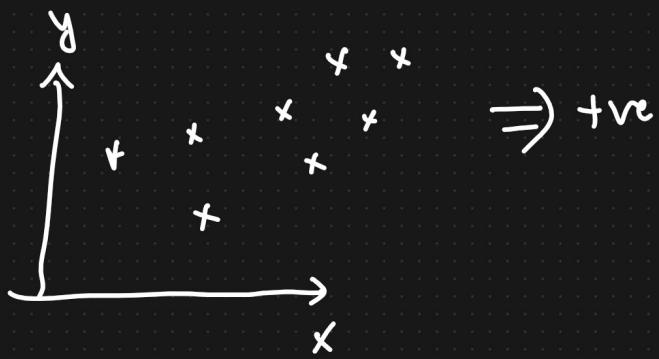
How much +ve } → Value  
How much -ve } =

## ② Pearson Rank Correlation Coefficient

$$\text{Cov}(x, y) = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{n-1}$$

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

$$\{ \text{Cov}(x, x) = \text{Var}(x) \}$$



No covariance

220

$$\boxed{-1 \rightarrow 1}$$

## Pearson Correlation Coefficient

$$f(x,y) = \frac{\text{Cov}(x,y)}{\sigma_x \sigma_y} \Rightarrow \begin{cases} -1 \text{ to } 1 \\ \uparrow \\ \downarrow \\ \text{Positive} \\ \text{Correlated} \\ \text{Unrelated} \\ \text{Negative} \end{cases}$$

$X$	$Y$
2	3
3	4
4	5
5	6

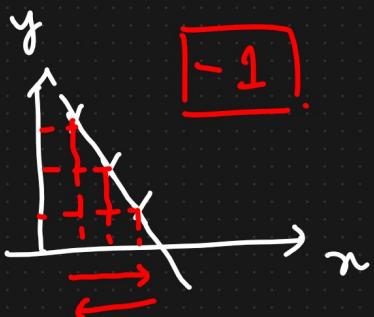
$$\bar{x} = 3.5 \quad \bar{y} = 4.5$$

$$\frac{\text{cov}(x, y)}{\sigma_x \sigma_y} = \boxed{-1 \rightarrow 1} \quad \checkmark$$

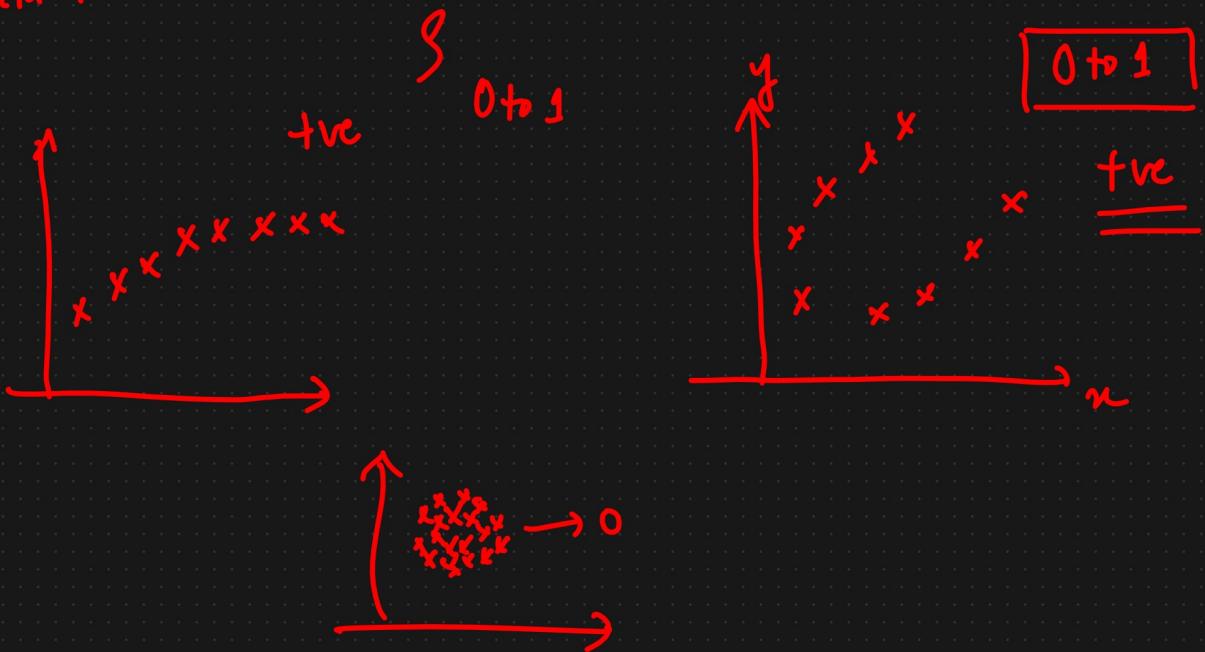
1

$$\begin{array}{|c|c|} \hline X & Y \\ \hline \end{array} \quad + 0.75 \quad \begin{array}{|c|c|} \hline X & Z \\ \hline \end{array} \quad + 0.8 \quad \checkmark$$

Perfectly  
negatively  
correlated



$x \uparrow y \downarrow$   
 $x \downarrow y \uparrow$



## Spearman Rank Correlation

$$= \frac{\text{Cov}(x, y)}{R\sigma_x R\sigma_y}$$

Feature Important



Domain Knowledge



X	Y	$R_x$	$R_y$
1	2	4	4
5	6	1	2
3	4	2	3
2	8	3	1

500 features

Feature Selection

① Independent features

② Dependent features

In  
Dependent  
features

## (2) Log Normal Distribution

No. of Study	No. of play	Time to sleep	O/P P/F
7	2	7	P
1	8	7	F

→	7	2	7	P
→	1	8	7	F

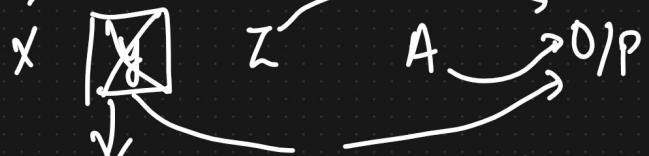
Mapped  $\rightarrow$  Run of the mouse

800 → 499 highly

Dependent

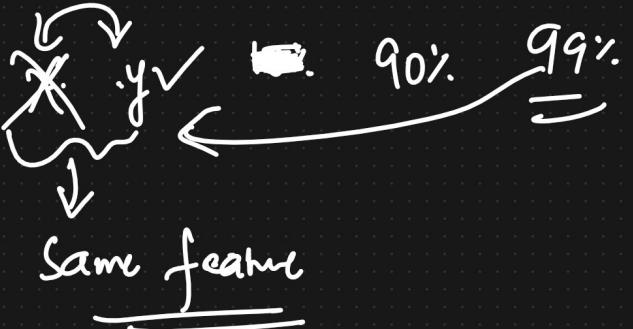
70%

$X \uparrow Y \downarrow$



Domain Expert

less →

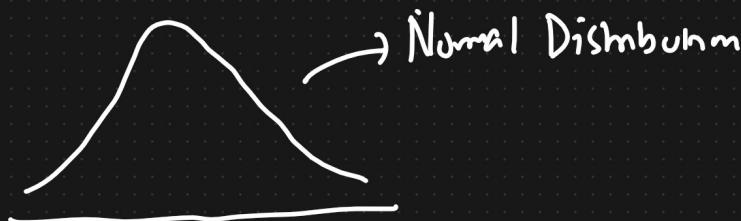


Same feature

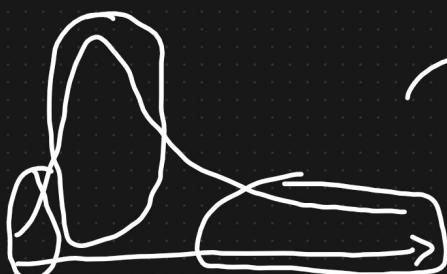
90% 99%

## \* Log Normal Distribution

$$\left\{ \begin{array}{c} \downarrow \\ 68 \end{array} \right. \left\{ \begin{array}{c} \downarrow \\ 95 \end{array} \right. \left\{ \begin{array}{c} \downarrow \\ 99.7 \end{array} \right. \}$$



Normal Distribution



Log Normal Distribution

Length of the  
comment



$$\Rightarrow y_{us} \Rightarrow \log(x)$$

Feature Transformation

Normality distributed data.

$$X \sim \text{Log Normal Distribution } (\mu, \sigma)$$



$$\text{Normal Distribution } \ln(x)$$

Some people

long comment  
Small comment  
Medium

Eg: V People writing comments In }  
My channel

Eg: Wealth distribution

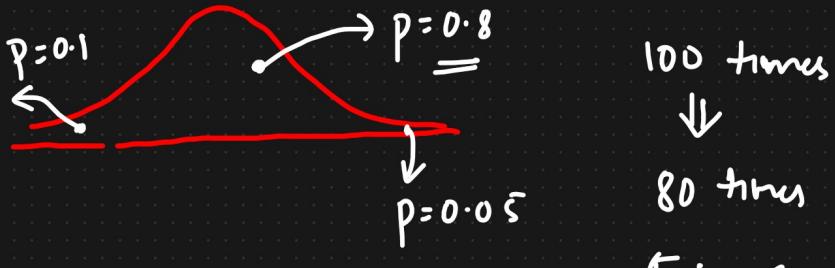
# Inferential Statistics

## ① P Value {probability value}

↓ Eg



Space bar



## Hypothesis Testing

Tossing a fair coin

100 times

$$\Pr(H) = \frac{1}{2}$$

$$\Pr(T) = \frac{1}{2}$$



50 times Head ✓

80 times Tail

$(H_0)$   $\xrightarrow{\text{True}}$  Null Hypothesis  $\xrightarrow{\text{True}}$  Treats Equals or Same  
 $\xrightarrow{\text{The coin is fair}}$  }  $\alpha = 1 - 0.95$

0.05 Alternate Hypothesis : The coin is not fair

$$= 0.05$$

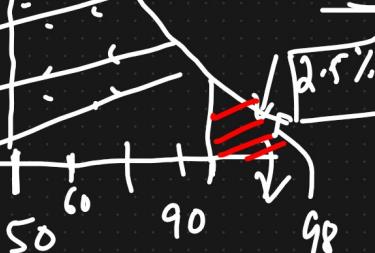
Confidence Interval = 95%

95% Symmetrical

Experiment

$$P < 0.05$$

$$= 0.05$$



9 Heads

$$\alpha = 0.35$$

$$2.5\%$$

40

$$\alpha = 0.05$$

50

60

70

80

90

98

↓ Domain Expert

$0.65$   
=

Reject the Null Hypothesis }  
Accept the Alternative Hypothesis }  
No. it is not a  
fair coin

Eg: Person has committed a crime

Null Hyp ( $H_0$ ) = He is innocent

Altan ( $H_1$ ) = He is not innocent

↳ Experiment  $\rightarrow$  Collecting evidence  $\rightarrow$  Finger print, DNA,

↳ Hypothesis Testing  $\chi^2 = \boxed{\quad} \rightarrow \underline{\text{Judge}}$

$F_0 \Rightarrow$  is not innocent

Fair  
Rolling a Dice

$$\boxed{1 \\ 6}$$

$$100 \\ \downarrow \\ 16 \quad 16 \quad 16 \quad 16$$

Inferential Stats

$$\chi^2 = \boxed{0.1} \quad \underline{\underline{I}}$$

Sample data  $\xrightarrow{\text{Assumption}}$  Population Data.

Experiment  
hypothesis

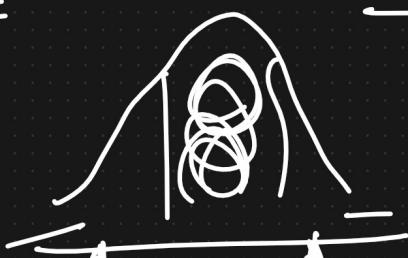
Average size of all the sharks in the world ??

C.I

90% C.I

$$\alpha = 10 \Rightarrow \underline{0.1}$$

|



$\mu < 0.1$

$$1 - 0.95 \Rightarrow \underline{0.05}$$



$$C.I = 85$$

$$\alpha = 0.15$$

