

statistics

$$P(\text{soccer}) = \frac{\text{soccer}}{\text{total}} = \frac{3}{10} \approx 0.3$$

event

↓

Sample space

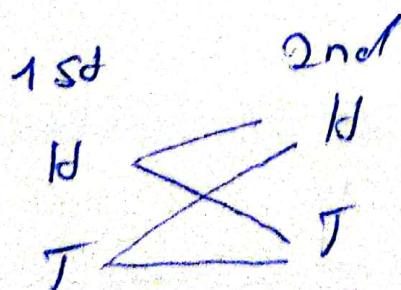
Coin example 1

on Experiment

probability of landing on heads : $P(\text{heads})$
 " " " " " " " " " tails : $P(\text{tails})$

$$P(\text{heads}) = \frac{1}{2} = 0.5$$

Coin exp 2 (2 coins)



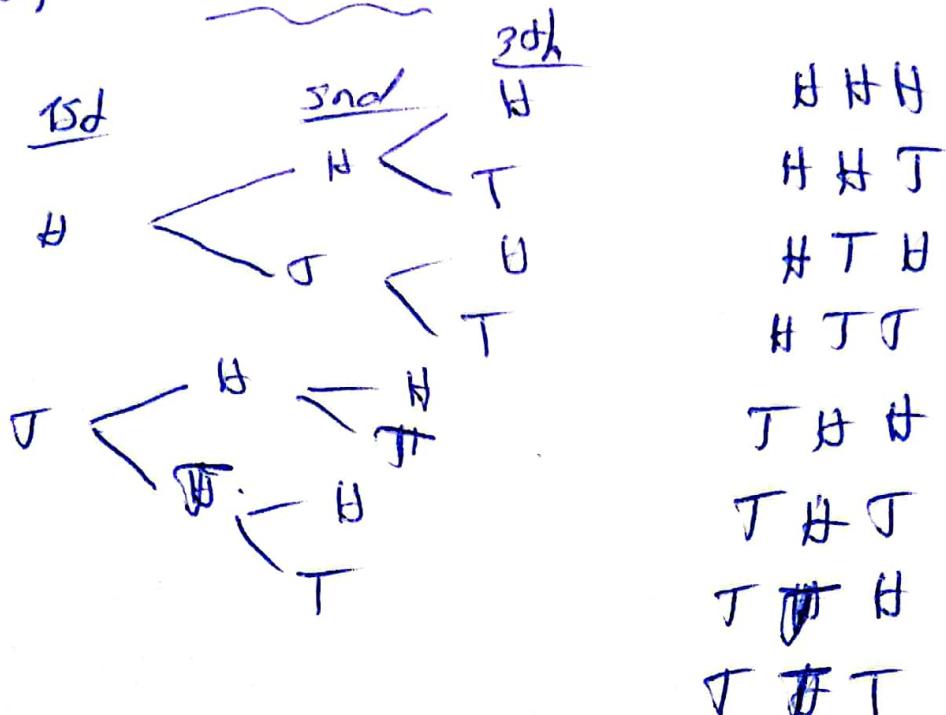
$$\text{both heads} = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} = 0.25$$

$$P(HH)$$

~~P(HHH)~~

Q

3 coins



$$P(HHH) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8} = 12.5\% = 0.125$$

Dice examples

$$P(6) = \frac{1}{6}$$

$$P(6+6) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

complement of probability

$$P(1-\text{soccer}) = 1 - \frac{3}{10} = \frac{7}{10} = 0.7$$

complement rule

$$P(A') = 1 - P(A)$$

(3)

$$P(HHH) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{8}$$

$$P(HHH') = 1 - \frac{1}{8} = \frac{7}{8}$$

$$P(\text{not } \sigma) = 1 - \frac{1}{8} = \frac{7}{8}$$

Sum of probability ~~(independent)~~

Soccer $P(S) = 0.3$

basketball $P(B) = 0.4$

$$P(S \text{ or } B) = 0.3 + 0.4 = 0.7$$

$$\overline{P(S \text{ or } B)} = P(S) + P(B)$$

$$\overline{P(S \cup B)} = P(S) + P(B) \rightarrow \boxed{\text{joint events}}$$

Dice exp:

what is the probability of obtaining an even number or a 5?

even: 2, 4, 6

$$P(5 \text{ or even}) = \frac{4}{6} = \frac{2}{3} \quad \frac{1}{6} + \frac{3}{6} = \frac{4}{6} = \boxed{\frac{2}{3}}$$

2 Dice

$A \cup B = \text{union}$

(4)

Sum of 7 or Sum of 10

$$1+6=7$$

$$2+5=7$$

$$3+4=7$$

$$4+3=7$$

$$5+2=7$$

$$6+1=7$$

$$\left. \begin{array}{l} 1+5=7 \\ 2+4=7 \\ 3+3=7 \\ 4+2=7 \\ 5+1=7 \\ 6+0=7 \end{array} \right\} \quad \begin{array}{l} 4+6=10 \\ 5+5=10 \\ 6+4=10 \end{array}$$

$$P(\text{Sum } 7 \text{ or Sum } 10) = \frac{6}{36} + \frac{3}{36} = \frac{9}{36} = \frac{1}{4}$$

difference of 2 or difference of 1

$$\left. \begin{array}{l} 6-5=1 \quad \textcircled{1} \\ 5-4=1 \quad \textcircled{2} \\ 4-3=1 \quad \textcircled{3} \\ 3-2=1 \quad \textcircled{4} \\ 2-1=1 \quad \textcircled{5} \end{array} \right\} \quad \begin{array}{l} 6-4=2 \quad 4-2=2 \\ 5-3=2 \quad 3-1=2 \\ 4-1=3 \quad 2-1=1 \end{array}$$

$$P(\text{d1 } \cup \text{ d2}) = \frac{10}{36} + \frac{8}{36} - \frac{12}{36} = \frac{7}{36} = \frac{1}{2} = 0.5$$

(disjoint events)

$$P(S) = 0.5$$

$$P(B) = 0.5$$

$$P(S \cup B) = P(S) + P(B) - P(S \cap B)$$

(5)

$$P(S) = 6$$

$$P(B) = 5$$

$$P(S \cap B) = 3$$

$$\text{total} = 10$$

$$P(S \cup B) = P(S) + P(B) - P(S \cap B)$$

$$= \frac{6}{10} + \frac{5}{10} - \frac{3}{10} = \frac{8}{10} = \frac{4}{5} = 0.8$$

Dice

sum \geq or diff 1

$$P(\text{sum } \geq \text{ or diff 1}) = \cancel{\frac{1}{6}} + \frac{6}{36} + \frac{10}{36} - \frac{2}{36} = \frac{12}{36}$$

Independence

chess moves \rightarrow dependence

osm \rightarrow ~~dependence~~ independence

~~total~~ Total = 100

~~S = 50~~

$$P(A \cap B) = P(A) \times P(B)$$

coin exps

what is the probability of landing on heads five times?

$$P(5 \text{ heads}) = \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} \times \frac{1}{2} = \frac{1}{32}$$

Dice exps

2 Dice both 6:

$$P(2 \text{ sixes}) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

$$P(10 \text{ sixes}) = \left(\frac{1}{6}\right)^{10}$$

Birthday Problem

probabilty of
no two people have the
same birthday

no two people have the
same birthday

50 people: 0.03
↓
no two people have the
same BD

23 people

number of
people

conditional probability

product rule (Independent events)

$$P(A \cap B) = P(A) \cdot P(B)$$

product rule (dependent events)

2 dices first 6 and sum 10

$$P(6 \cap \text{sum}10) = \frac{1}{6} \times \frac{1}{6} = \frac{1}{36}$$

↙ ↘

6 4

$$P(\text{sum}=10 | \text{1st } 6) = P(\text{1st } 6) \cdot P(\text{sum}=10 | \text{1st } 6)$$

$$P(A \cap B) = P(A) \cdot P(B|A)$$

↑
General product rule

notes

Total = 100

S = 50

room1 = 50

room2 = 50

$$P(S) = 0.4$$

(8)

$$P(\text{not } S) = 0.6$$

$$P(R|S) = 0.8$$

$$P(S \cap R) = P(S) \cdot P(R|S) = 0.4 \times 0.8 = 0.32$$

$$P(\text{not } S \cap R) = P(\text{not } S) \cdot P(R|\text{not } S)$$

$$= 0.6 \times 0.2 = 0.12$$

$$P(R|S) = 0.8$$

$$\begin{array}{c} P(S) = 0.4 \\ \swarrow \quad \searrow \\ R \qquad \text{not } R \end{array}$$

$$P(\text{not } R|S) = 0.2$$

$$\begin{array}{c} P(S) = 0.4 \\ \swarrow \quad \searrow \\ R \qquad \text{not } S \end{array} \quad P(R|\text{not } S) = 0.5$$

$$P(\text{not } R|\text{not } S) = 0.5$$

Bayes Theorem

$$P(\text{sick} | \text{diagnosed sick}) = \frac{\text{sick and diagnosed sick}}{\text{sick and diagnosed sick} + \text{healthy and diagnosed sick.}}$$

(9)

$$P(\text{sick}) = 0.01\%$$

sick = A

$$P(\text{not sick}) = 99.99\%$$

diagnosed sick = B

$$P(\text{diagnosed sick} | \text{sick}) = 99\%$$

$$P(\text{diagnosed sick} | \text{not sick}) = 1\%$$

$$P(A|B) = \frac{P(\text{sick}) \cdot P(\text{diagnosed sick} | \text{sick})}{P(\text{sick}) \cdot P(\text{diagnosed sick} | \text{sick}) + P(\text{not sick}) \cdot P(\text{diagnosed sick} | \text{not sick})}$$

$$P(A|B) = \frac{0.01 \times 0.99}{0.01 \times 0.99 + 0.9999 \times 0.1} = 0.0098$$

$$P(A) = 0.01\%$$

$$P(A') = 99.99\%$$

$$P(B|A) = 99\%$$

$$P(B|A') = 1\%$$

$$P(A|B) = \frac{P(A) \cdot P(B|A)}{P(A) \cdot P(B|A) + P(A') \cdot P(B|A')}$$

prior

$$P(A)$$

event

$$E$$

Posterior

$$P(A|E)$$

(10)

Naive assumption (in Bayes theorem)

$$P(A | w_1, \dots, w_n) = \frac{P(A) \cdot P(w_1 | A) \cdots P(w_n | A)}{P(A) \cdot P(w_1 | A) \cdots P(w_n | A) + P(A') \cdot P(w_1 | A') \cdots P(w_n | A')}$$

Exps:

1) Image recognitions

$$P(\text{Cat} | \text{image}) = P(\text{Cat} | \text{pixel}_1, \text{pixel}_2, \dots, \text{pixel}_n)$$

2) classification in medicine

$$P(\text{healthy} | \text{symptoms and history})$$

Random Variables

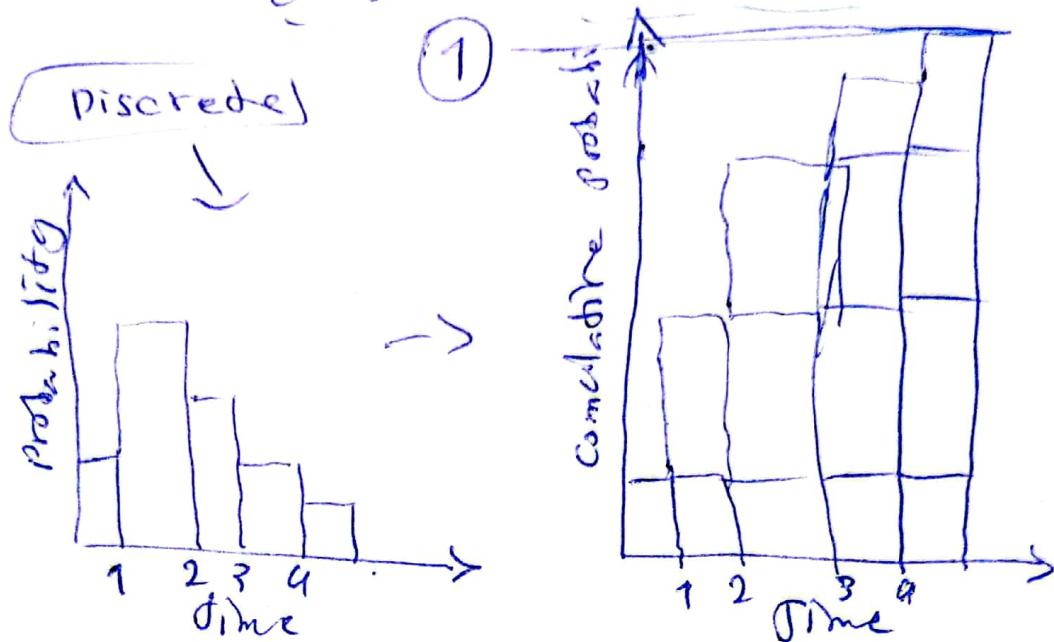
Exp, ~~flip a coin~~ $\begin{cases} H, P(H) = 0.5 \\ T, P(T) = 0.5 \end{cases}$

if $H \rightarrow X = 1$

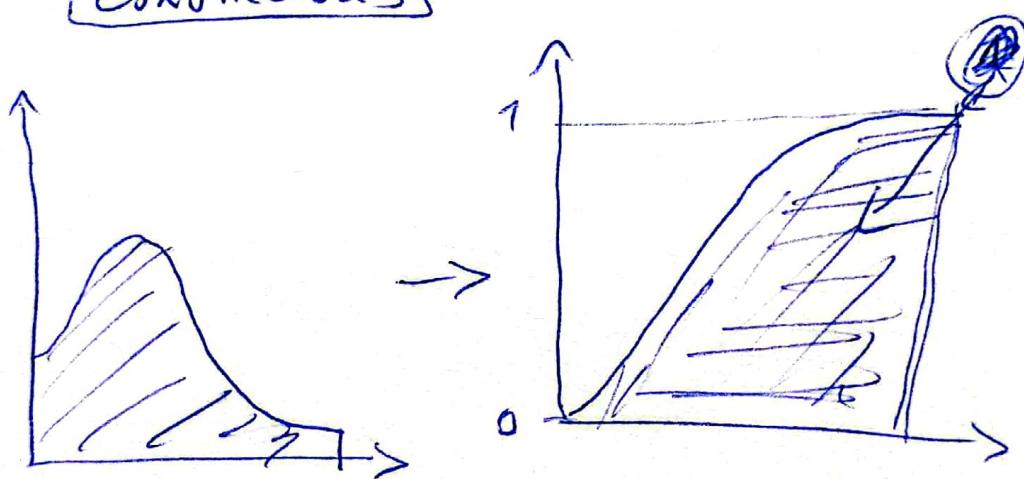
if not $H \rightarrow X = 0$

X is a random variable
and in this case can be $\underline{1}$ or $\underline{0}$

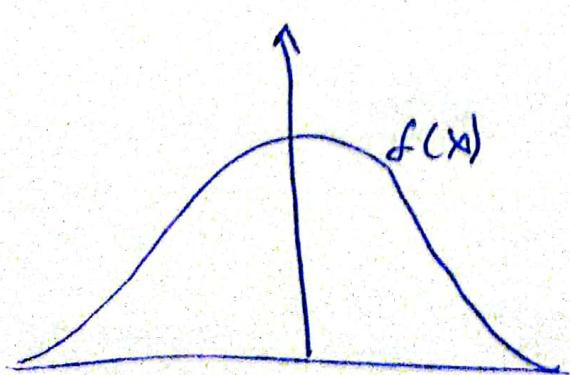
Cumulative distribution



Continuous

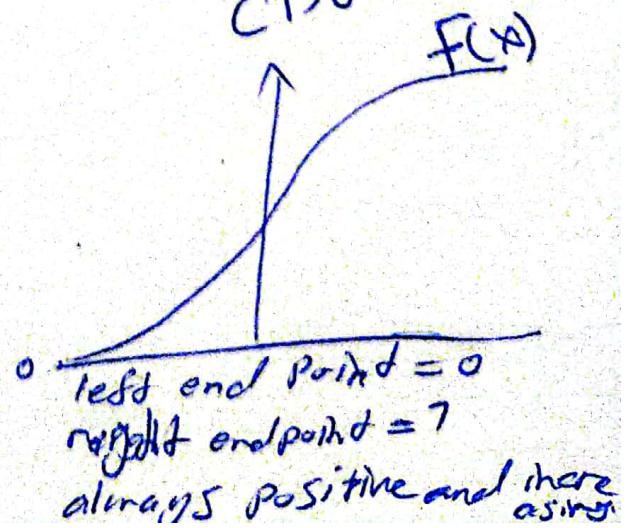


PDF



area = 1
always positive

CDF



"Binomial distribution"

Binomial coefficient $\binom{n}{k}$ counts all the combinations for landing k heads in n coin tosses.

$$\binom{n}{k} = \binom{n}{n-k}$$

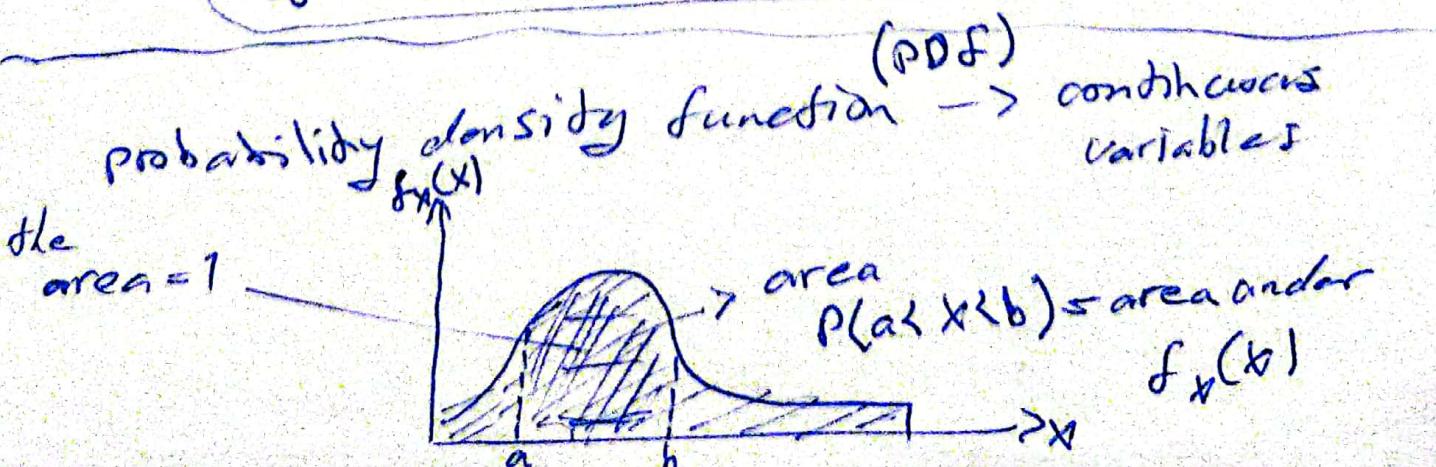
$$\binom{n}{k} = \binom{n}{n-k} = \frac{n!}{k!(n-k)!}$$

~~PDF~~
probability mass function \rightarrow discrete variables

$$P_X(x) = \binom{n}{x} p^x (1-p)^{n-x}, \quad x=0,1,2,\dots,n$$

$X \sim \text{Binomial}(n, p)$

n and p are called the parameters of the binomial distribution.



Random variables

Discrete (can take only a countable number of values)

Continuous (take values on an interval)

Random Variables eg:

~~X=2~~

$X = \text{number of defective items in a shipment}$

↙
uncertain outcome

Deterministic Variables

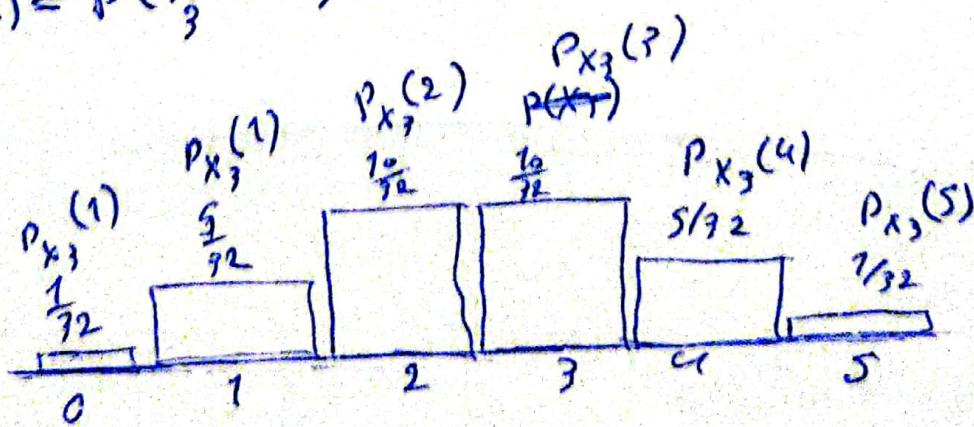
$X=2 \quad f(X)=X^2$



fixed outcome

probability mass function (PMF)

$$P_X(x) = P(X=x), \quad x \in \{0, 1, 2, 3, 4, 5\}$$



Binomial distribution with very large n = normal distribution.

$$\text{normal (Gaussian) distribution} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

$$X \sim N(\mu, \sigma^2)$$

X is normal distributed

$$\text{Standardization} \Rightarrow Z = \frac{X - \mu}{\sigma}$$

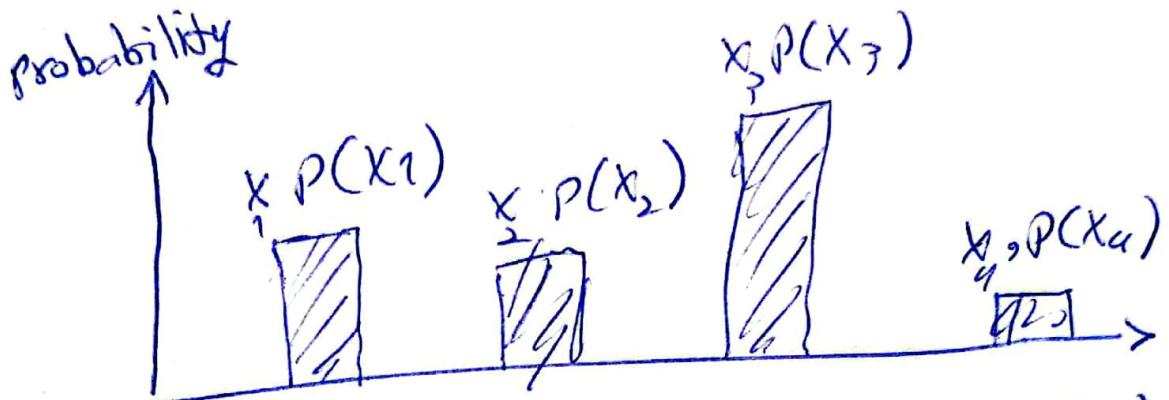
"Joint Distribution (Discrete)"

~~$P_{xy}(x,y) = P(X=x, Y=y)$~~

Expected value of a scenario

(15)

$$E[X] = x_1 p(x_1) + x_2 p(x_2) + x_3 p(x_3) + \dots + x_n p(x_n)$$



$$E[f(x)] = f(x_1)p(x_1) + f(x_2)p(x_2) + f(x_3)p(x_3) + \dots + f(x_n)p(x_n)$$

$$\text{Exp: } E[X^2] = \frac{1+4+9+16+25+36}{6} = \frac{91}{6} = \cancel{\frac{E[X^2]}{6}}$$

$$\underline{E[X] = E[X_i] + E[X_j]}$$

$$\text{in general: } E[X_1 + X_2] = E[X_1] + E[X_2]$$

$$\underline{E[\text{matches}] = E[X_1] + E[X_2] + \dots + E[X_n] = 7}$$

(10)

Variance for molar

$$\begin{aligned}\text{Var}(X) &= E[(X-\mu)^2] \\ &= E[X^2] - E[X]^2\end{aligned}$$

Standard deviations:

$$\text{std}(X) = \sqrt{\text{Var}(X)}$$

$$= \sqrt{E[(X-\mu)^2]}$$

$$= \sqrt{E[X^2] - E[X]^2}$$

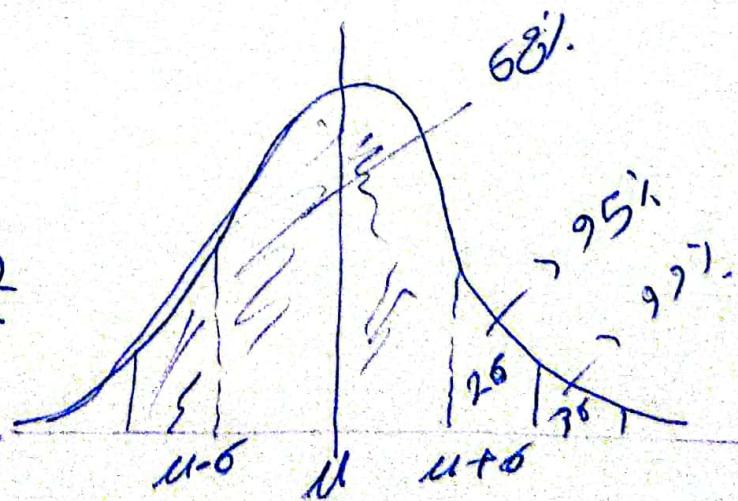
Normal distribution

μ = center of the bell

σ = spread of the bell

$\rightarrow X \sim N(\mu, \sigma^2)$

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2}}$$



"Sum of Gaussians"

Sum of two or more Gaussian distributions
is also Gaussian distributed.

$$R = (J+L) \sim N\left(\frac{\mu_J + \mu_L}{2}, \frac{\sigma_J^2 + \sigma_L^2}{2}\right)$$

In General: $W = aX + bY$

independent

$$\begin{cases} X \sim N(\mu_X, \sigma_X^2) \\ Y \sim N(\mu_Y, \sigma_Y^2) \end{cases}$$

$$W \sim N\left(a\mu_X + b\mu_Y, a^2\sigma_X^2 + b^2\sigma_Y^2\right)$$

"Moments of a distribution"

$$E[X] = P_1 X_1 + P_2 X_2 + \dots + P_n X_n$$

$$E[X^2] = P_1 X_1^2 + P_2 X_2^2 + \dots + P_n X_n^2$$

$$E[X^3] = P_1 X_1^3 + P_2 X_2^3 + \dots + P_n X_n^3$$

$$\dots$$

$$E[X^n] = P_1 X_1^n + P_2 X_2^n + \dots + P_n X_n^n$$

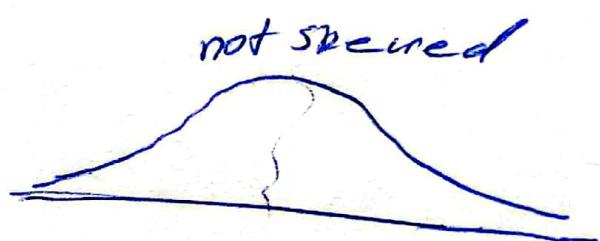
SKEWNESS

$$\text{Skewness} = E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right]$$

$$E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] > 0$$



$$E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] = 0$$



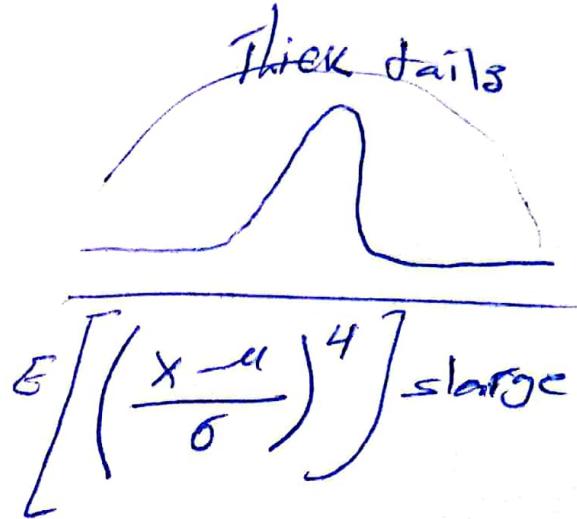
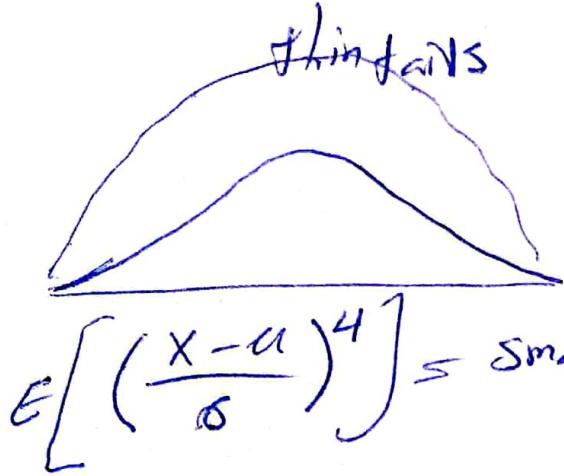
$$E \left[\left(\frac{X - \mu}{\sigma} \right)^3 \right] < 0$$



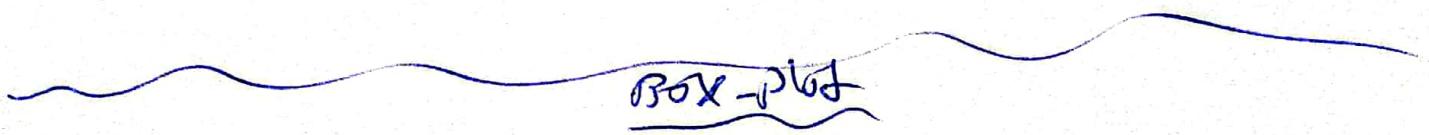
Kurtosis

~~$E[X^4]$~~

$$\text{kurtosis} = E \left[\left(\frac{X - \mu}{\sigma} \right)^4 \right]$$



↓
cattosis



$$Q_1 = g_{0.25}$$

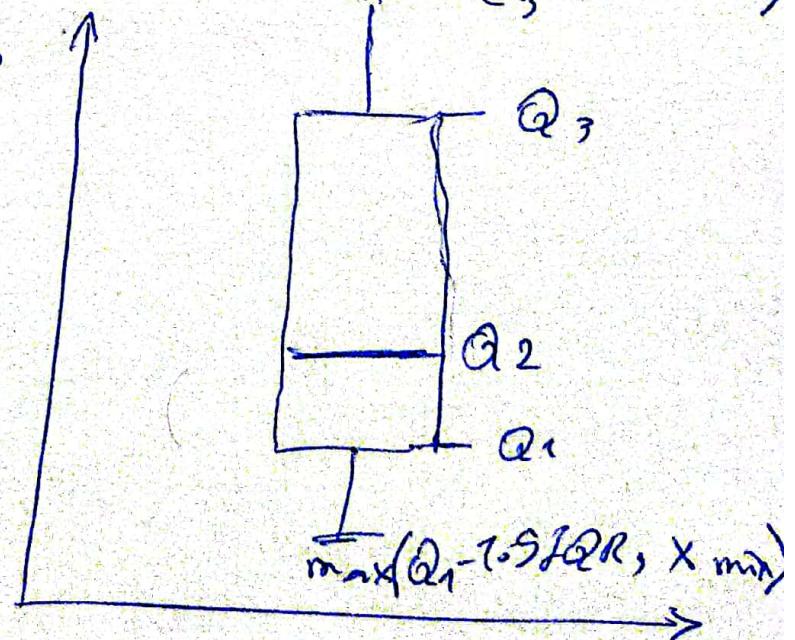
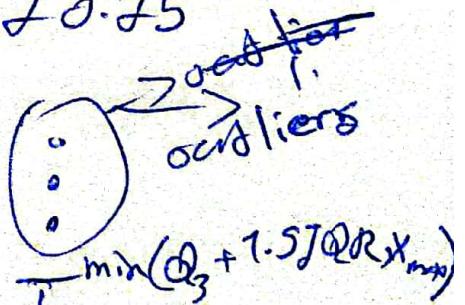
$$Q_2 = g_{0.5}$$

$$Q_3 = g_{0.75}$$

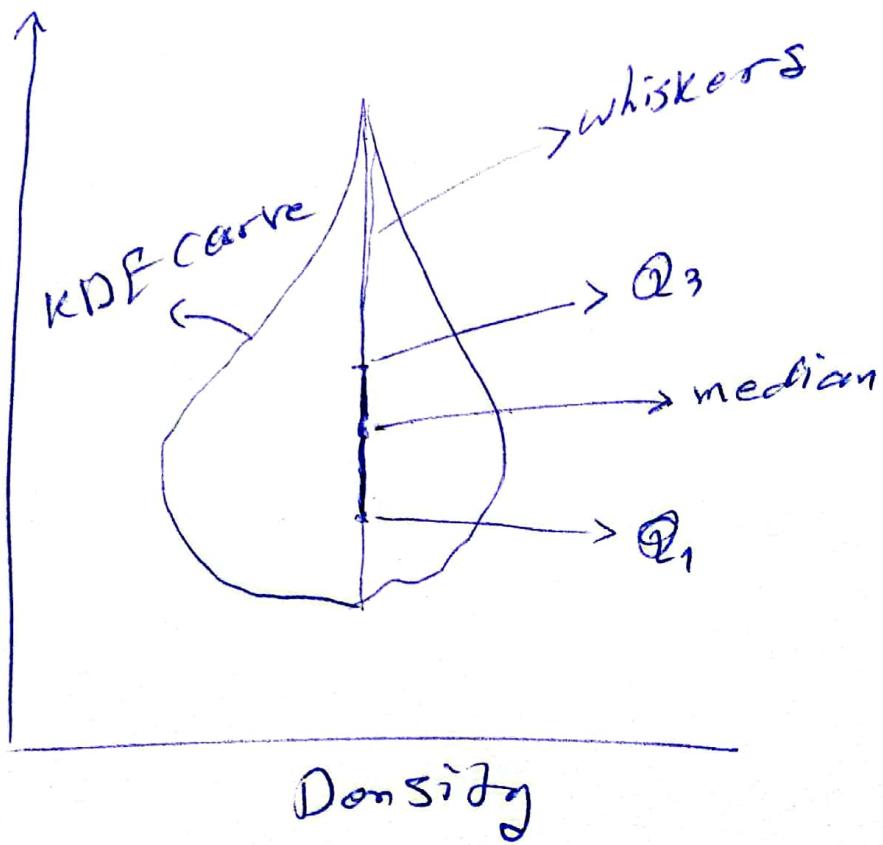
$$\text{Interquartile range (IQR)} = Q_3 - Q_1$$

$$x_{\min} = ?$$

$$x_{\max} = ?$$

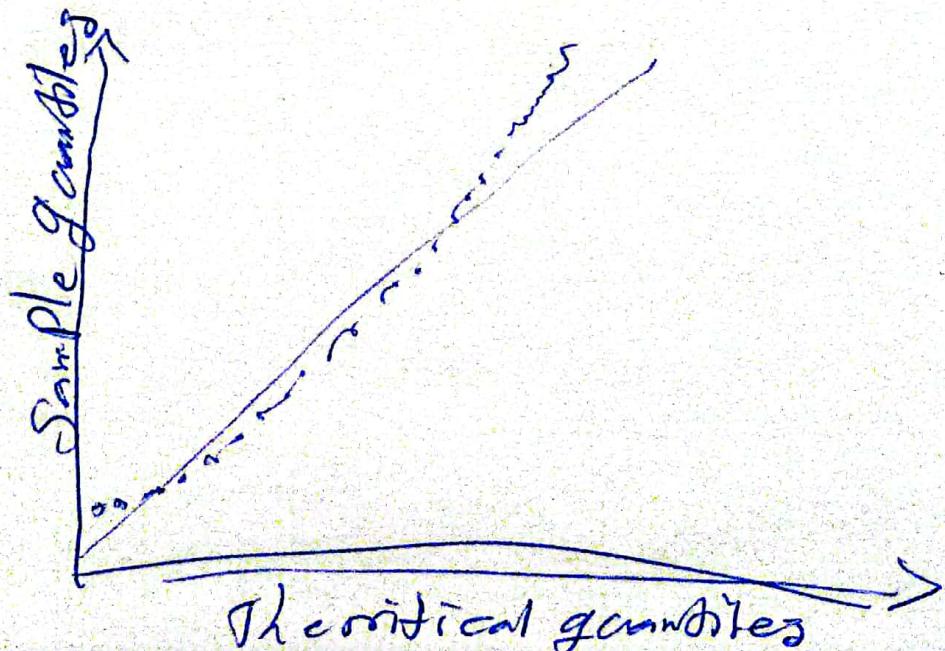


Violin-Plot



QQ plots

check normality (Gaussian distribution)



"joint distribution" (discrete)

(27)

$$P_{xy}(x,y) = P(X=x, Y=y)$$

for dependent discrete
variables

Exps:

		Age (Y) $n=10$			
		45	46	47	48
height (X) $n=10$	7	$2/10$	$2/10$	0	0
	8	0	0	$2/10$	$1/10$
	9	0	0	0	$3/10$

(all combination of probabilities)
(discrete joint)

$$P_{xy}(x,y) = P(X=x, Y=y) = P(X) \cdot P(Y)$$

independent random variables
discrete

marginal distribution

$$P_y(y_i) = \sum_i P_{xy}(x_i, y_i)$$

$$P_x(x_i) = \sum_j P_{xy}(x_i, y_j)$$

"discrete Conditional distribution"

$$P_{Y|X=i}(y) = P(Y=y | X=i)$$

→ joint PDF of X and Y

$$P_{Y|X=x}(y) = \frac{P_{xy}(x, y)}{P_x(x)}$$

↓

conditional PDF of Y

→ ~~PDF of X~~

→ marginal distribution of X

"Conditional conditional distribution" (23)

conditional PDF of Y

$$f_{Y|X=x}(y) = \frac{f_{XY}(x,y)}{f_X(x)}$$

marginal distribution
of X

joint PDF of X and Y

COVariance

$$\text{Cov}(X, y) = \sum xy$$

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

positive covariance

negative covariance

"Covariance of a probability distribution" (24)

$$\text{Var}(X) = \sum_{i=1}^N (x_i - \bar{x})^2 \cdot P(x_i)$$

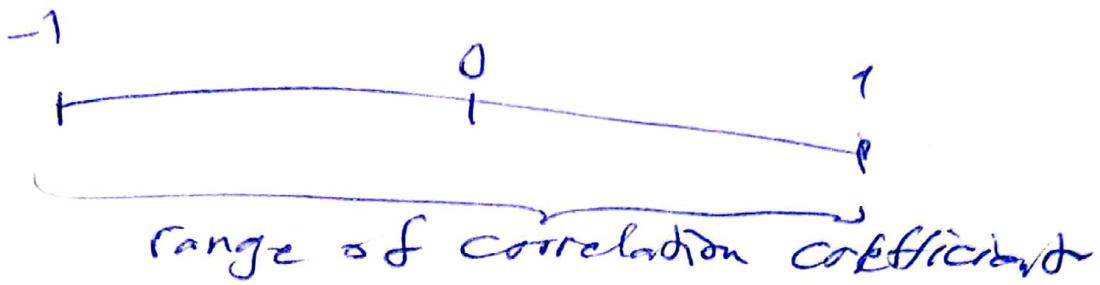
$$\begin{aligned}\text{Cov}(X, Y) &= \sum P_{XY}(x_i, y_i)(x_i - \bar{x})(y_i - \bar{y}) \\ &= E[XY] - E[X]E[Y]\end{aligned}$$

Covariance matrix

$$\Sigma = \begin{bmatrix} X & Y \\ \downarrow & \downarrow \end{bmatrix} \begin{bmatrix} \text{Var}(X) & \text{cov}(X, Y) \\ \text{cov}(X, Y) & \text{Var}(Y) \end{bmatrix}$$

Sigma or covariance matrix

Correlation Coefficient



$$\text{cov}(X, Y)$$

$$\text{Correlation coefficient} = \frac{\text{cov}(X, Y)}{\sigma_X \cdot \sigma_Y}$$

it is actually the standardized covariance.

or

$$\text{cor}(X, Y)$$

$$\text{Corr-coeff} = \frac{\text{cov}(X, Y)}{\sqrt{\text{Var}(X)} \cdot \sqrt{\text{Var}(Y)}}$$

Multivariate Gaussian Distribution

For a single variable

$$f_X(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

~~25~~
if \mathbf{W}, \mathbf{H} were independent

$$f_{HW}(H, w) = f_H(H)f_w(w)$$

$$= \frac{1}{\sqrt{2\pi\sigma_H^2}} e^{-\frac{1}{2} \frac{(h-\mu_H)^2}{\sigma_H^2}} \cdot \frac{1}{\sqrt{2\pi\sigma_w^2}} e^{-\frac{1}{2} \frac{(w-\mu_w)^2}{\sigma_w^2}}$$

$$= \frac{1}{2\pi\sigma_H\sigma_w} e^{-\frac{1}{2} \left(\frac{(h-\mu_H)^2}{\sigma_H^2} + \frac{(w-\mu_w)^2}{\sigma_w^2} \right)}$$



$$= \frac{1}{2\pi \det \Sigma} \exp \left(-\frac{1}{2} \left(\begin{bmatrix} h \\ w \end{bmatrix} - \boldsymbol{\mu} \right)^T \Sigma^{-1} \left(\begin{bmatrix} h \\ w \end{bmatrix} - \boldsymbol{\mu} \right) \right)$$

$$f_X(x_1, x_2, \dots, x_n) = \frac{1}{(2\pi)^{n/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-\boldsymbol{\mu})^T \Sigma^{-1} (x-\boldsymbol{\mu})}$$

↓
 mean vector
 $\boldsymbol{\mu} = [u_1, u_2, \dots, u_n]^T$

covariance matrix spread
 of the bell

$|\Sigma| \rightarrow$ determinant of
the cov matrix

"Statistics"

Every dataset you work with in machine learning
is a sample (not a population)

Population size = N

Sample size = n

Population mean = μ

Sample mean = \bar{X}

Population proportion

$P_s = \frac{\text{number of items with a given characteristic (}x\text{)}}{\text{population (}n\text{)}}$

Sample proportion

$$\hat{P} = \frac{x}{n}$$

$x \rightarrow$ number of items in sample
 \downarrow
sample size

Population Variance Sample Variance

$$\sigma^2 = \frac{1}{N} \sum (x - \mu)^2$$

↓
population mean

population size

Variance estimation

~~$\sigma^2 = \frac{1}{n} \sum (x - \mu)^2$~~

Sample(X) $\rightarrow \text{Var}(X) = \frac{\sum (x - \bar{x})^2}{n-1} \rightarrow$ estimated variance

~~SS (1)~~

Sample Variance

$$\text{Var}(X) = \frac{1}{n-1} \sum (x - \bar{x})^2$$

law of large numbers

As the sample size increases, the average of the sample will tend to get closer to the average of the entire population.

n = number of samples

X_i = some estimate X for a sample size i

as $n \rightarrow \infty$

$$\frac{1}{n} \sum_{i=1}^n X_i \rightarrow E[X] = \mu_X$$

conditions:

- 1) Sample is randomly drawn
- 2) Sample size must be sufficiently large.
- 3) Independent observations.

Central limit theorem

As the number of observation (n) increases, the probability distribution of the sample becomes closer to a Gaussian distribution, even if the original variables (population) are not normally distributed.

$$\bar{X} \approx \mu = np$$

$$\sigma^2 \approx np(1-p)$$

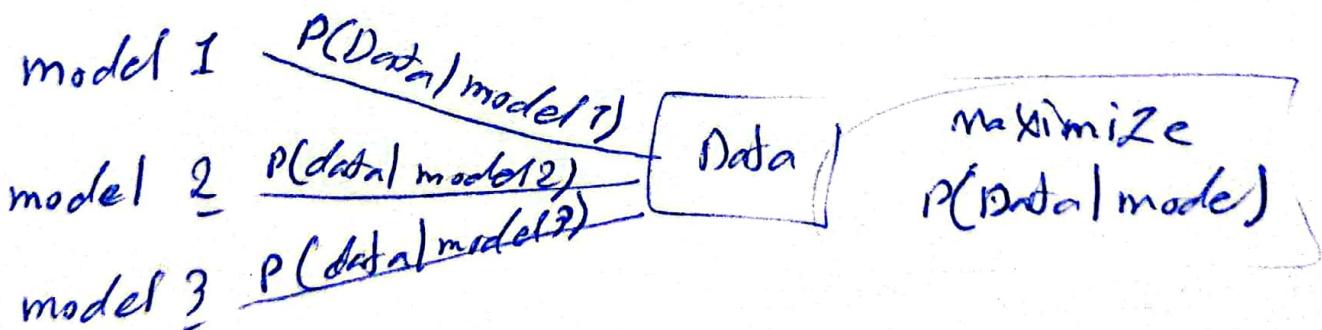
$$\text{As } n \rightarrow \infty \quad \frac{\frac{1}{n} \sum_{i=1}^n x_i - E[X]}{\sigma_X} \xrightarrow{} \sqrt{n} \sim N(0, 1^2)$$

$$\text{As } n \rightarrow \infty \quad \frac{\sum_{i=1}^n x_i - nE[X]}{\sqrt{n} \sigma_X} \xrightarrow{} N(0, 1^2)$$

"Point estimation"

"maximum likelihood estimation" (MLE)

"find the model that most likely produced the data"



Maximum Likelihood Bernoulli Example

Data: 8 heads 2 tails coin

Models:

Model 1: $P(H) = 0.7$	$P(T) = 0.3$	$\left. \begin{array}{l} \text{Model 2: } P(H) = 0.5 \\ P(T) = 0.5 \end{array} \right\} \begin{array}{l} \text{Model 3: } \\ P(H) = 0.3 \\ P(T) = 0.7 \end{array}$

→ maximum P(model 1 / Data)

an abetter model

$$\text{if } p = P(H) \quad \text{Likelihood } L(p; \text{Data}) = p^8(1-p)^2$$

$$= \log(p^8(1-p)^2) = \cancel{8\log(p)} + 2\log(1-p)$$

*↑
log likelihood*

we try to maximize the log likelihood, which
can be interpreted as maximizing the
likelihood.

$$\frac{d}{dp} f(8\log(p) + 2\log(1-p)) = \frac{8}{p} + \frac{2}{1-p} (-1)$$

$$\text{if } \frac{8}{p} + \frac{2}{1-p} (-1) = 0 \Rightarrow \hat{p} = \frac{8}{10}$$

if $n = \text{heads}$

$k = \text{heads}$

$$x_1 \cdot x_2 \cdot \dots \cdot x_{n-1} \cdot x_n$$

$$x = (x_1, \dots, x_n)$$

$$x_i \sim \text{Bernoulli}(p)$$

likelihood

$$L(p; x) = P_p(x=x) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$\text{if } x_i = 1, p^{[x_i]} (1-p)^{[1-x_i]} = p$$

$$\text{if } x_i = 0, p^{[x_i]} (1-p)^{[1-x_i]} = (1-p)$$

$$\sum_{i=1}^n x_i = \# \text{ heads}$$

$$n - \sum_{i=1}^n x_i = \# \text{ tails}$$

so:

$$L(p; x) = P_p(x=x) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$= p^{\left(\sum_{i=1}^n x_i\right)} (1-p)^{\left(n - \sum_{i=1}^n x_i\right)}$$

likelihood

$$CC(p; x) = \log \left(CP^{\sum_{i=1}^n x_i} \cdot (1-p)^{n-\sum_{i=1}^n x_i} \right) \quad (33)$$

$$= \left(\sum_{i=1}^n x_i \right) \cdot \log(p) + \left(n - \sum_{i=1}^n x_i \right) \cdot \log(1-p)$$

find the maximum:

$$\frac{d}{dp} L(p; x) = \frac{d}{dp} \left(\left(\sum_{i=1}^n x_i \right) \log(p) + \left(n - \sum_{i=1}^n x_i \right) \cdot \log(1-p) \right)$$

$$= \frac{\sum_{i=1}^n x_i}{p} + \frac{n - \sum_{i=1}^n x_i}{1-p} (-1) = 0$$

$$\Rightarrow \hat{p} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x} \rightarrow \text{mean of the population}$$

maximum likelihood: Gaussian Example:

The best model is the mode where the mean and variance of the population is equal to sample

mean and variance of the population

Suppose $X = (X_1, X_2, \dots, X_n)$ Gaussian distributed (34)

So, $X_i \stackrel{i.i.d}{\sim} N(\mu, \sigma^2)$

mean μ

Variance $= \sigma^2$

$$\begin{aligned} L(\mu, \sigma; X) &= \prod_{i=1}^n f_{X_i}(x_i) = \\ &= \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \frac{(x_i - \mu)^2}{\sigma^2}} \\ &= \frac{1}{(\sqrt{2\pi})^n \sigma^n} e^{-\frac{1}{2} \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2}} \end{aligned}$$

find the values of μ and σ that maximize the likelihood $L(\mu, \sigma; X)$

$$C(\mu, \sigma) = \log(L(\mu, \sigma; X))$$

$$L(\mu, \sigma) = \log \left(\frac{1}{(\sqrt{2\pi})^n \sigma^n} e^{-\frac{1}{2} \frac{\sum_{i=1}^n (x_i - \mu)^2}{\sigma^2}} \right)$$

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

(35)

Partial derivative of ~~μ~~ ^{log likelihood} with respect do μ and σ

~~$\frac{\partial}{\partial \mu}$~~

$$\frac{d}{d\mu} L(\mu, \sigma) = \frac{-1}{2} \frac{\sum_{i=1}^n 2(x_i - \mu)}{\sigma^2} (-1)$$

$$\left\{ \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - \sum_{i=1}^n \mu \right) = \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - n\mu \right) \right.$$

$$\frac{d}{d\sigma} L(\mu, \sigma) = -\frac{n}{\sigma} - \frac{1}{2} \left(\sum_{i=1}^n (x_i - \mu)^2 \right) (-2) \frac{1}{\sigma^3}$$

$$= -\frac{n}{\sigma} + \left(\sum_{i=1}^n (x_i - \mu)^2 \right) \frac{1}{\sigma^3}$$

equating the partial derivatives do ~~$\neq 0$~~ 0

$$\frac{d}{d\mu} L(\mu, \sigma) = \frac{1}{\sigma^2} \left(\sum_{i=1}^n x_i - n\mu \right) = 0$$

since $\sigma > 0$ $\sum_{i=1}^n x_i - n\mu = 0$

so:
$$\hat{\mu} = \frac{\sum_{i=1}^n x_i}{n} = \bar{x}$$

sample mean

LME for mean

$$\frac{d}{d\sigma} L(\mu, \sigma) = -\frac{n}{\sigma} + \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \frac{1}{\sigma^3} = 0$$

(50)

since $\sigma > 0$:

$$\frac{d}{d\sigma} L(\mu, \sigma) = -n + \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \frac{1}{\sigma^2} = 0$$

$$\hat{\mu} = \bar{x} \rightarrow$$

$$\frac{d}{d\sigma} L(\mu, \sigma) = -n + \left(\sum_{i=1}^n (x_i - \bar{x})^2 \right) \frac{1}{\sigma^2} = 0$$

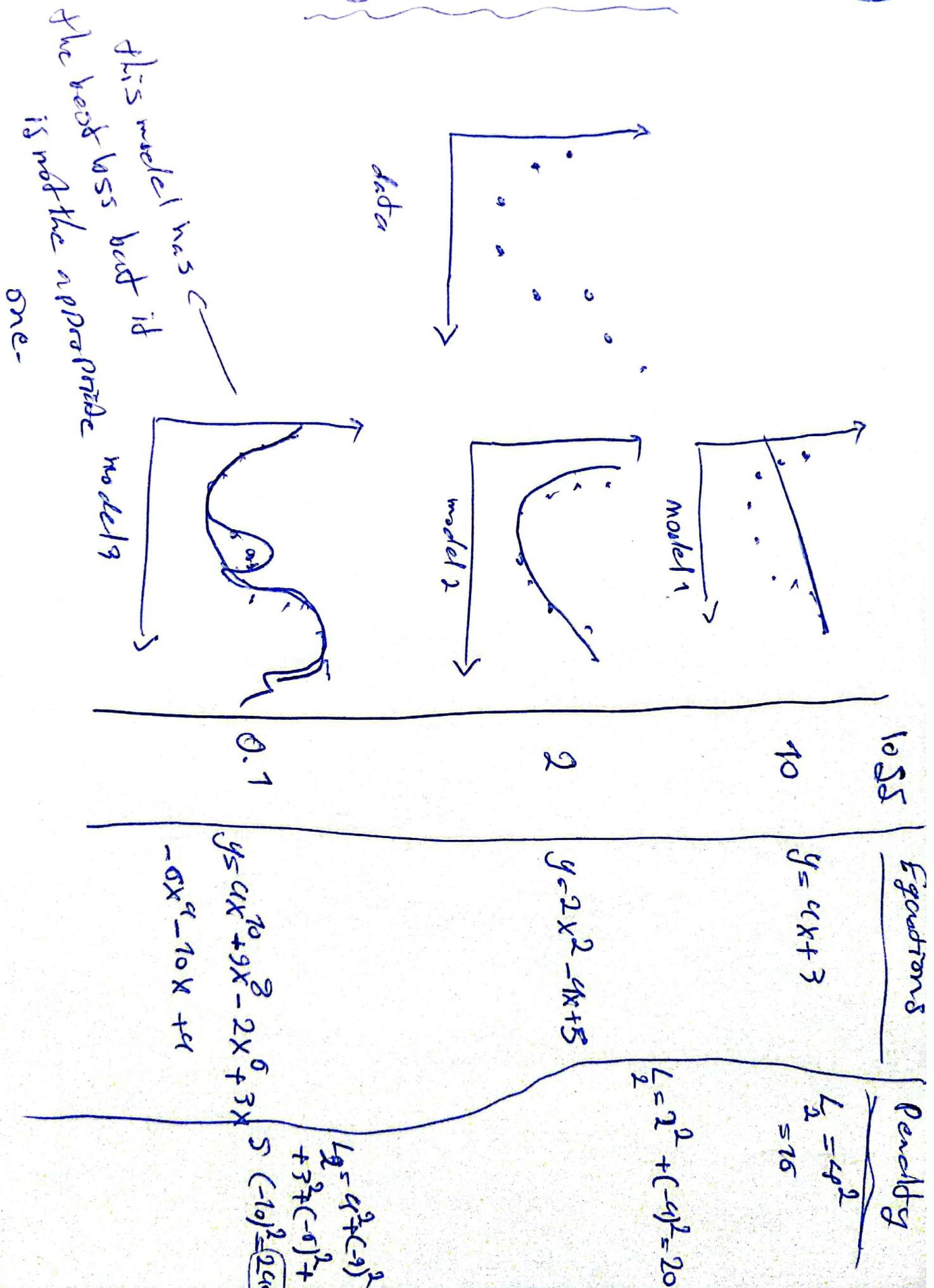
$$\sigma^2 = \frac{\sum (x_i - \bar{x})^2}{n}$$

50)

$$\sigma = \sqrt{\frac{\sum (x_i - \bar{x})^2}{n}}$$

LME for the standard deviation

Polynomial regression



38

near loss
= loss + penalty

start

$$10 + 10 = 20$$

$$2 + 20 = 22$$

best model to

describe the system

Regularization term in general

$$\text{model: } y = a_n x^n + a_{n-1} x^{n-1} + \dots + a_1 x + a_0$$

$$\log \text{loss} = CL$$

$$\text{Error} = a_n^2 + a_{n-1}^2 + \dots + a_1^2$$

L₂ regularization term: a small number

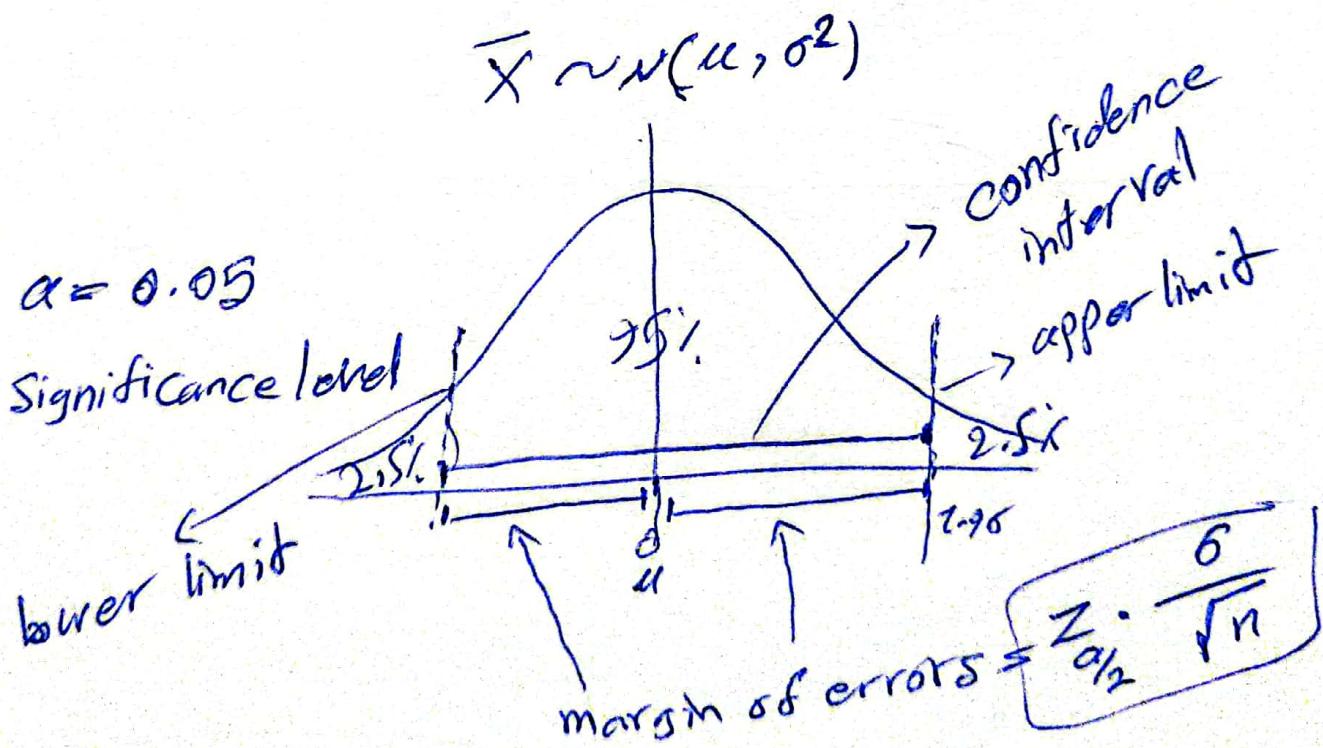
Regularization parameter: 1

Statistical Inference

Bayesian

frequentist
(maximum likelihood...)

Z-distribution



if $\mu=0$ and $\sigma=\tau$ ($X \sim N(0, \tau)$)

then $I_{\alpha/2} = 1.96$
critical value

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \rightarrow \text{Standard error}$$

(40)

confidence interval calculation

1) find the sample means
→ pop size

$N=6000$ adults

$$\begin{cases} n = 49 \\ \bar{x} = 170 \text{ cm} \end{cases} \quad \sigma \rightarrow \text{sample std}$$

$$\sigma = 25 \text{ cm}$$

2) find the desired confidence level ($1-\alpha$)

find 95% confidence interval

3) get the critical value ($Z_{\alpha/2}$)

$$\text{for 95\%} \rightarrow Z_{\alpha/2} = 1.96$$

4) find the standard error $\left(\frac{\sigma}{\sqrt{n}} \right)$

$$\frac{\sigma}{\sqrt{n}} = \frac{25 \text{ cm}}{\sqrt{49}} = \frac{25}{7}$$

5) find the margin of error

$$Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} = 1.96 \times \frac{25}{7} = 7$$

margin of error

6) Add/subtract the margin of error to the sample mean

$$\text{confidence intervals } \bar{x} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

mean of the population
→ pop size

$$\text{confidence interval} = 170 \text{ cm} \pm 7 \text{ cm}$$

$$50 \text{ } 163 \text{ cm } 177 \text{ cm}$$

* - The confidence interval contains the true population parameter approximately 95% of the time.

\nearrow for means
confidence interval - student's t-distribution

when σ known

$$\text{Confidence interval} = \bar{X} \pm Z_{\alpha/2} \cdot \frac{\sigma}{\sqrt{n}} \rightarrow \text{num sample}$$

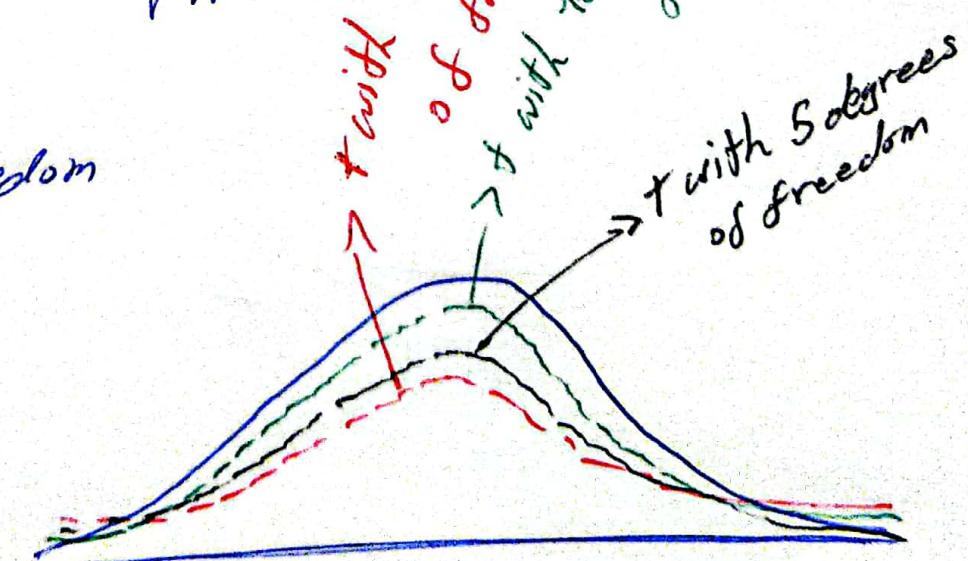
when σ unknown

$$\text{Confidence interval} = \bar{X} \pm t_{\alpha/2} \cdot \frac{s}{\sqrt{n}} \rightarrow \text{sample-std}$$

+ with 1 degrees of freedom
+ with 10 degrees of freedom

degrees of freedom

$$n-1$$



Confidence interval for proportion (Q2)

$\checkmark X = 24$ $\hat{p} = \frac{X}{n} = \frac{24}{30} = 80\%$

$n = 30$

confidence interval = $\hat{p} \pm \text{margin of error}$
 for proportion

$$\text{margin of error} = I_{\alpha/2} \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}$$

Standard error

margin of error = $1.96 \times \sqrt{\frac{0.8 \times 0.2}{30}}$

= 0.14

confidence interval = 0.8 ± 0.14

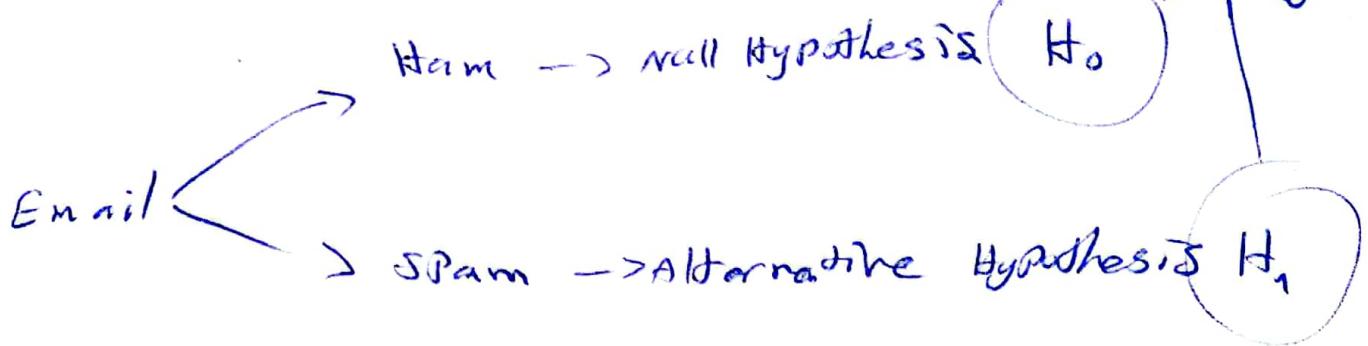
$0.66 < p < 0.94$

60% < P < 94%

Hypothesis Testing

43

Expt:



Email is assumed to be Ham

Type I and Type II errors

Expt:

false positive

if you send a regular email to spam Box: Type I error

if you send a spam Email to the regular Box: Type II error

false negative

decision	H_0 True	H_0 False
reject H_0	Type I error	correct
Don't reject H_0	correct	Type II error

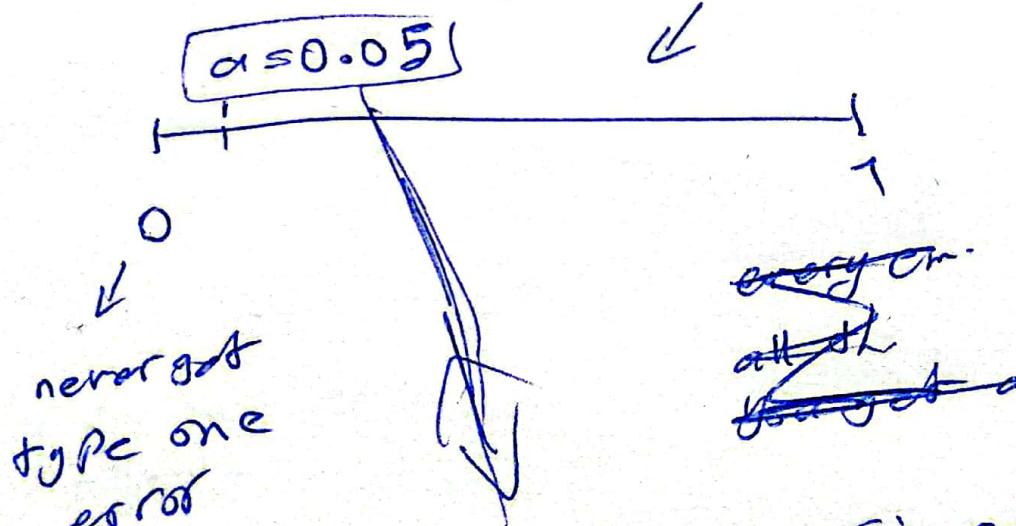
(c4)

Type one error is more Dangerous

Ex: Sending a regular email to spam is worse than not sending a spam email to the regular inbox.

- what is the greatest probability of type I error you willing to tolerate?

↓
significance level (α)



a Ham email is spam ~~on~~ 5% of times

Right-left- and two-tailed test 95

Right-Tailed Test $\rightarrow H_0: \mu = \mu_0$ vs. $H_1: \mu > \mu_0$

Left-Tailed Test $\rightarrow H_0: \mu = \mu_0$ vs. $H_1: \mu < \mu_0$

Two-Tailed Test $\rightarrow H_0: \mu = \mu_0$ vs. $\underline{H_1: \mu \neq \mu_0}$

P-value

A P-value is the probability, assuming H_0 is true, that the test statistic takes on a value as extreme or more extreme than the value observed.

if $p\text{-value} < \alpha$ reject the H_0 (and accept H_1)

if $p\text{-value} > \alpha$ do not reject H_0

$T(X)$: test statistic T : observed statistic

$$H_0: \mu = \mu_0$$

