

Probability: A measure of how likely an event is to occur

$\star$  Sum of Probabilities

$$\begin{cases} \text{Disjoint events: } P(A \cup B) = P(A) + P(B) & \text{Mutually exclusive} \\ \text{Joint events: } P(A \cup B) = P(A) + P(B) - P(A \cap B) & \text{Non-mutually exclusive} \end{cases}$$

Independence: When the occurrence of one event does not affect the probability of the occurrence of another event.

Product Rule for Independent events:  $P(A \cap B) = P(A) \cdot P(B)$

Conditional Probability: Calculating the probability of an event happening given that another event has already happened.

e.g:

$$P(HH \mid 1st \text{ is } H) = \frac{1}{2} \quad P(HH \mid 1st \text{ is } T) = \frac{0}{2}$$

The general Product Rule:

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

\*When independent:  $P(B|A) = P(B)$

Bayes Theorem: We have 1 mil people in the population and the illness only affects

- 1 of every  $10^5$  people. The test that the doctor gives you is 99% effective  
(out of every 100 sick people, 99 would be diagnosed sick and 1 would be diagnosed healthy)  
and out of 100 healthy people, 99 would be healthy and 1 would be sick)

You went to the doc and you tested sick! Are you really sick or not?



$$P(\text{sick} | \text{diagnosed sick}) = \frac{\text{sick and diagnosed sick}}{\text{diagnosed sick}} = \frac{99}{99 + 9999} = 0.0098$$

$\rightarrow A = \text{sick}$   $B = \text{diagnosed sick}$

Bayes formula

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

$\underbrace{P(B|A) + P(A \cap B) - P(B)}$

$$P(B|A) = 99\% \quad P(A) = 0.01\% \quad P(A') = 99.99\% \quad P(B|A') = 1\%$$

$$\rightarrow P(A|B) = 0.0098$$

Prior: The original probability that you can calculate not knowing anything else  
 $P(A)$

Event: Gives us information about the probability  $E$

Posterior:  $P(A|E)$

e.g.: Prior =  $P(\text{sick})$  event = diagnosed positive Posterior =  $P(\text{sick} | \text{diag pos})$

\* Posterior is a much more accurate probability

Subject: \_\_\_\_\_

Date: \_\_\_\_\_

e.g.2: Prior:  $P(\text{spam})$  event: Email contains 'Lottery' Posterior:  $P(\text{spam} \mid \text{lottery})$

e.g.3: Prior:  $P(\text{spam})$  event: Email contains 'winning' Posterior:  $P(\text{spam} \mid \text{winning})$

How to combine these two events?

Naive assumption (foundation of naive bayes): The appearances of lottery and winning are independent (even though they are not!)

$$P(\text{spam} \mid \text{lottery \& winning}) = \frac{P(\text{spam}) \cdot P(\text{lottery \& winning} \mid \text{spam})}{P(\text{spam}) \cdot P(\text{lottery \& winning} \mid \text{spam}) + P(\text{ham}) P(\text{lottery \& winning} \mid \text{ham})}$$

$$\frac{P(A \cap B) = P(A)P(B)}{P(\text{spam}) P(\text{lottery} \mid \text{spam}) P(\text{winning} \mid \text{spam}) + P(\text{ham}) P(\text{lottery} \mid \text{ham}) P(\text{winning} \mid \text{ham})}$$

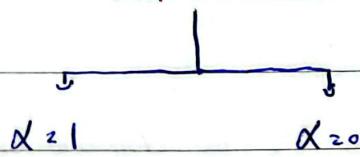
in general (naive bayes):

$$P(A_1 | B_1, B_2, \dots, B_n) = \frac{P(A) P(B_1 | A) \dots P(B_n | A)}{P(A) P(B_1 | A) \dots P(B_n | A) + P(\bar{A}) P(B_1 | \bar{A}) \dots P(B_n | \bar{A})}$$

Random Variable: Variables that can take many values → e.g.: Temperature

$X$  = number of heads

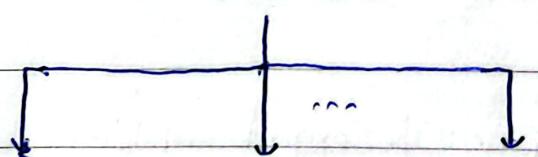
Flip a coin



$$P(x=1) = 0.5$$

$$P(x=0) = 0.5$$

Flip 10 coins



$$P(x=10) = 0.5^{10}$$

$$P(x=9) = ?$$

$$P(x=0) = 0.5^{10}$$

\*  $\sum P(x=i) = 1$

Custom random vars: Define  $X$ :

$$P(x=1) = 0.5 \quad P(x=-7) = 0.2$$

$$P(x=\pi) = 0.3$$

## Discrete & Continuous Random Vars:

Discrete: Can take only a countable number of values

\* can be finite: Number of heads when flipping a coin

\* can be infinite:  $X = \text{Number of } \underset{\text{coin}}{\text{flips}} \text{ until we get our first head}$

Continuous: Can take values on an interval, e.g.: Wait time until the next bus arrives  $\rightarrow 1 \text{ min? } 1.01 \text{ min? } 1.0001 \text{ min? } \dots$

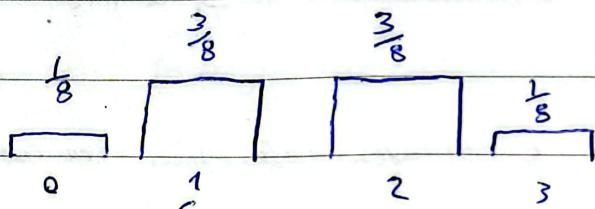
## Random Vars vs Deterministic Vars

Deterministic: Fixed outcome. e.g.:  $X = 2$  or  $P_{Xn} = n^2$  (once it's defined, it's like that forever)

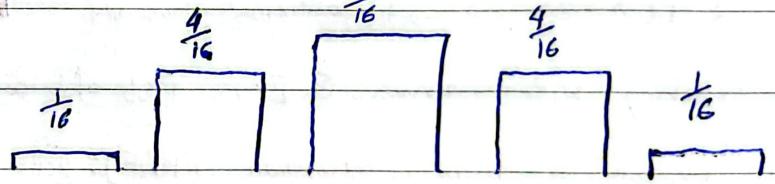
Random: Uncertain outcome. e.g.:  $X = \text{number of defective items in a shipment}$

## Probability Distributions (Discrete)

$X_1$ : number of heads in 3 coin tosses



$X_2$ : number of heads in 4 coin tosses



$$P(X_2 = n) \quad n = 0, 1, 2, 3, 4 \quad \xrightarrow{=} P_X(n) \rightarrow (\text{PMF: Probability mass function})$$

\*  $P_X(n) \geq 0$       \*  $\sum_n P_X(n) = 1$

\* PMF is for Discrete Vars

Subject: \_\_\_\_\_

Date: \_\_\_\_\_

**Binomial Distribution**: Number of heads in  $n$  coin tosses?  $P(H) = P$

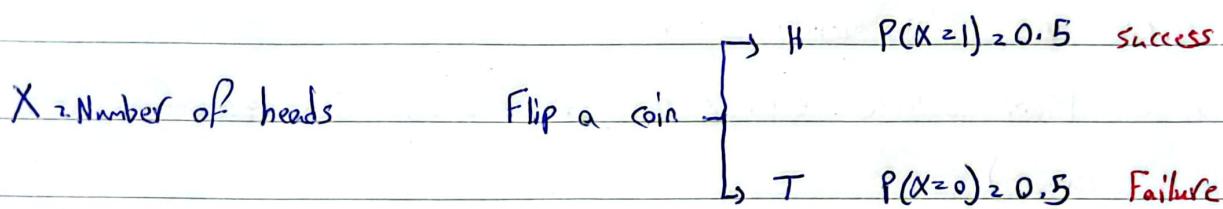
Event:  $X \sim n \rightarrow n$  heads in  $n$  tosses

$$P_X(n) = \binom{n}{n} P^n (1-P)^{n-n}, n=0,1,2,\dots \quad X \sim \text{Binomial}(n, P)$$

Ex: What is the probability of getting three ones when rolling a dice 5 times?

$$P(\text{one}) = \frac{1}{6} \quad P(\text{not one}) = \frac{5}{6} \quad P_X(3) = \left(\frac{1}{6}\right)^3 \left(\frac{5}{6}\right)^2 = 0.032$$

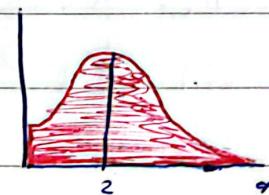
**Bernoulli distribution**:  $X \sim \text{Bernoulli}(P)$   $P$ : Probability of success



### Probability Distribution (Continuous)

\* Discrete  $\rightarrow$  a list. Continuous  $\rightarrow$  an interval

\* Red area (area under the curve)  $\approx 1$



\*  $P(X=2) = 0 \rightarrow$  In continuous, we are

working with intervals. Exact Prob of a given point is zero.

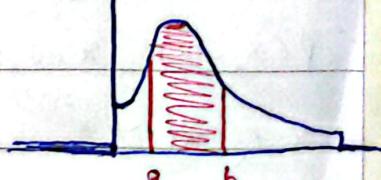
**Probability Density Function (PDF)**: Tells you the rate you accumulate Prob around each point,  $P(a < X < b) = \text{area under } f_X(x)$

$f_X(x)$  needs to satisfy: 1. Defined for all numbers (it can be zero)

$$2 - f_X(x) \geq 0$$

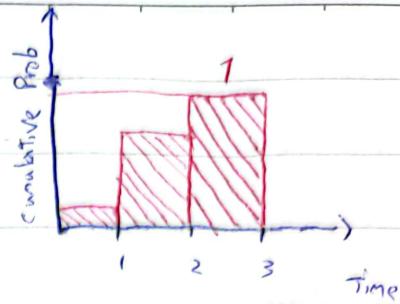
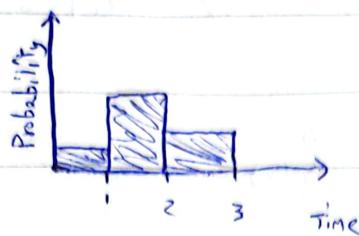
$$3 - \text{Area under } f_X(x) \approx 1$$

\* PDF only defined for continuous vars



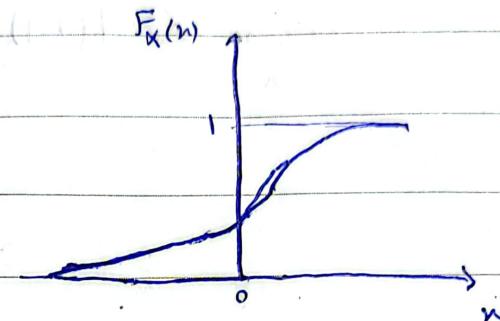
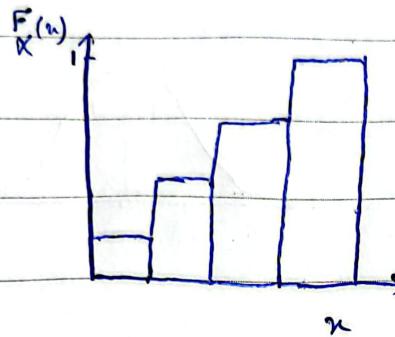
## Cumulative Distribution

\* Starts at 0 and ends at 1



\* The CDF shows how much prob the variable has accumulated until a certain value

$$CDF(n) = P(X \leq n) \rightarrow \text{Defined for every real number}$$



Properties:

- 1 -  $0 \leq F_X(n) \leq 1$
- 2 - Left Point = 0
- 3 - Right Point = 1
- 4 - Never Decreases

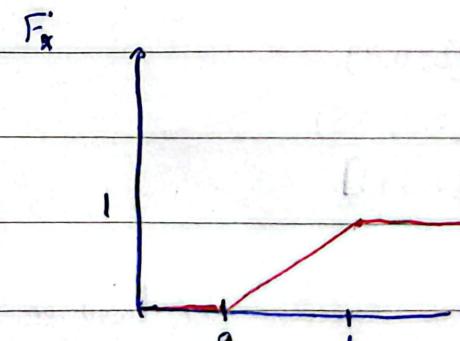
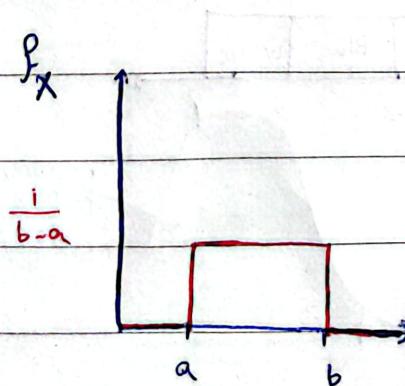
**Uniform Distribution:** A continuous random var can be modeled with a uniform distribution if all possible values lie in an interval and have same frequency of occurrence

Parameters:

a: beginning of the interval

b: end

$$f_X(x) = \begin{cases} \frac{1}{b-a} & a < x < b \\ 0 & \text{ow} \end{cases}$$



$$F_X(n) = \begin{cases} 0 & x < a \\ \frac{x-a}{b-a} & a \leq x < b \\ 1 & b \leq x \end{cases}$$

Normal (Gaussian) Distribution:

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

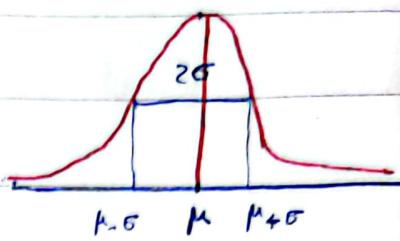
Parameters:

$\mu$ : center of the bell

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

Scaling constant

$\sigma$ : spread of  $x$

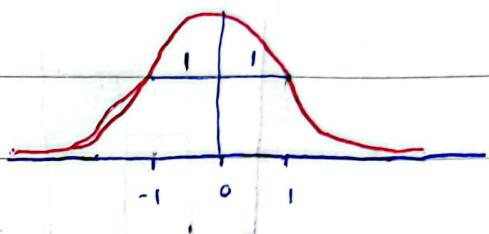


\* Range: All numbers

\* Symmetrical

Standard Normal Distribution:  $X \sim N(0, 1)$

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

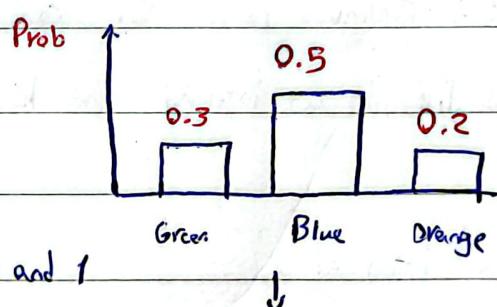


Standardization: Convert to std-Normal dist

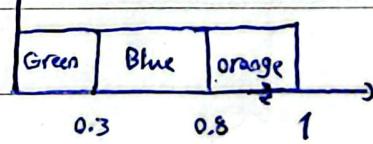
e.g.:  $X$  distributes normally with  $\mu = 2$ ,  $\sigma = 2.5 \rightarrow Z = \frac{X-\mu}{\sigma} = \frac{X-2}{2.5}$

Sampling from a dist:

Create a random sample of data that follows a dist



1. Generate a random number between 0 and 1



2. Find out which interval the number belongs to:

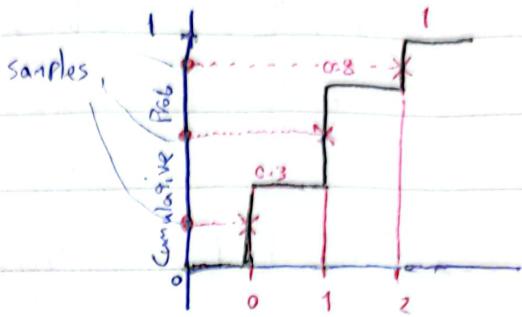
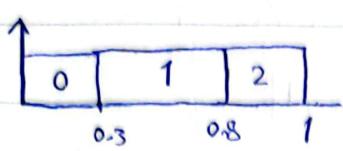
- $[0, 0.3)$
- $[0.3, 0.8)$
- $[0.8, 1]$

3. Assign the correct color based on the interval.

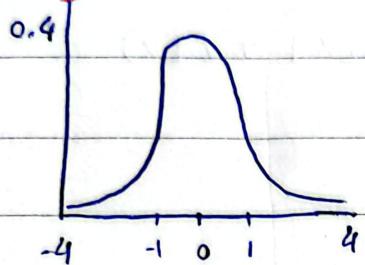
Date:

Subject:

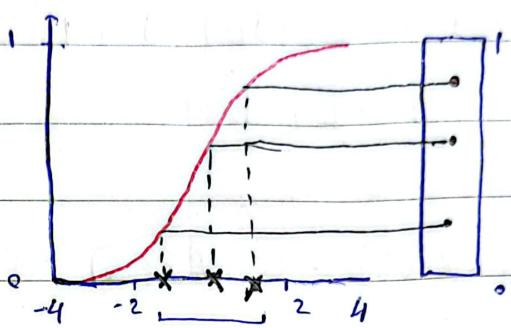
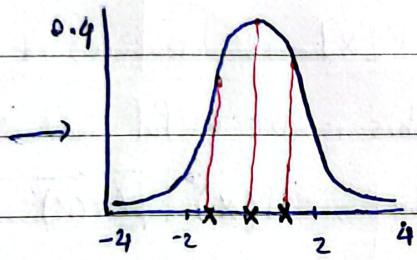
Another way:



Sampling from a normal dist



CDF

Now we have  
these Points

★ Mean = Avg Median = data in the middle (in terms of ordered position)

Mode = the value with highest prob (we can have more than one mode  $\rightarrow$  Multimodal)★  $E[X] = \text{Mean/Balancing point of a distribution} \rightarrow \text{Expected Value}$ ★  $E[X] = \sum_i P(x_i)x_i$  (expected value of a var)       $E[f(x)] = \sum_i f(x_i)P(x_i)$   
(expected value of a function)★  $E[aX+b] = aE[X]+b$        $E[b] = b$ Sum of Expectation:  $E[x_1 + x_2 + \dots + x_n] = E[x_1] + E[x_2] + \dots + E[x_n]$ 

e.g.: Imagine that you have 8 billion unique names in a bag and we have 8 billion people around the world. We are going to travel the world and give each person one piece of paper. What is the expected number of correct assignments?

Subject:

Date:

N People ( $N=8$  billion)

$$E[X_{\text{Mashes}}], E[X_{\text{Person 1}}] + E[X_{\text{Person 2}}] + \dots + E[X_{\text{Person } N}]$$

$$= \frac{1}{n} + \frac{1}{n} + \dots + \frac{1}{n} = \boxed{1}$$

Variance :  $E[X] = \mu$ ,  $\text{Var}(X) = E[(X - \mu)^2]$

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

Properties of the Variance

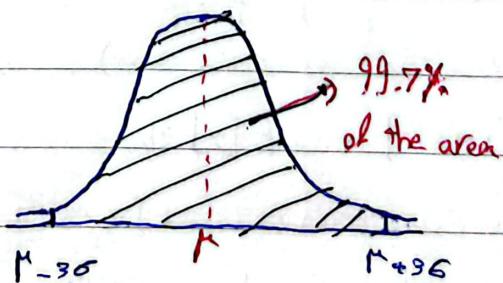
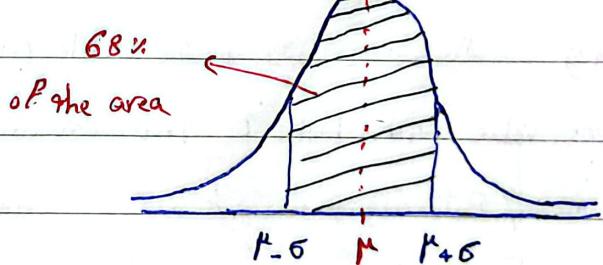
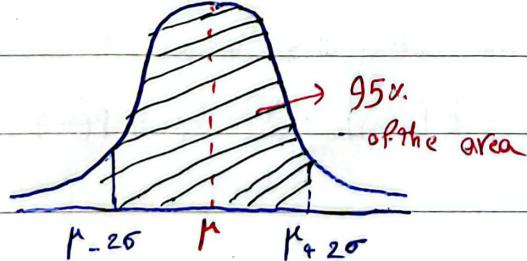
\*  $\text{Var}(aX + b) = a^2 \text{Var}(X)$

**Standard Deviation** : Say  $X$  is measured in unit  $Y$ ,  $E[X]$  is measured in  $Y$ ,  $\text{Var}(X)$  is measured in  $Y^2$  → but sometimes (sometimes) this is not useful and we want to have measurements with same unit as variable →  $\text{Std}(X) = \sqrt{\text{Var}(X)}$

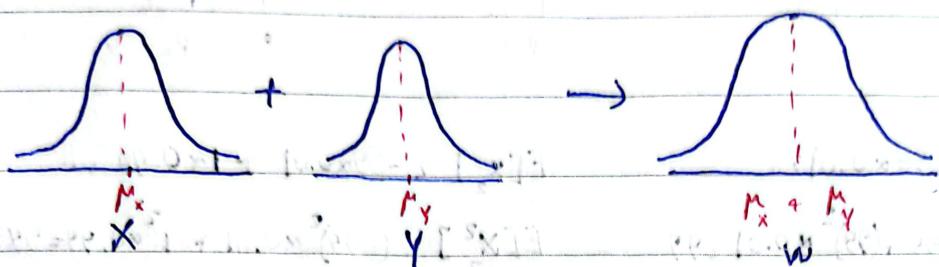
Normal distribution: 68-95-99.7 Rule

$\mu$ : Expected value

$\sigma$ : Std



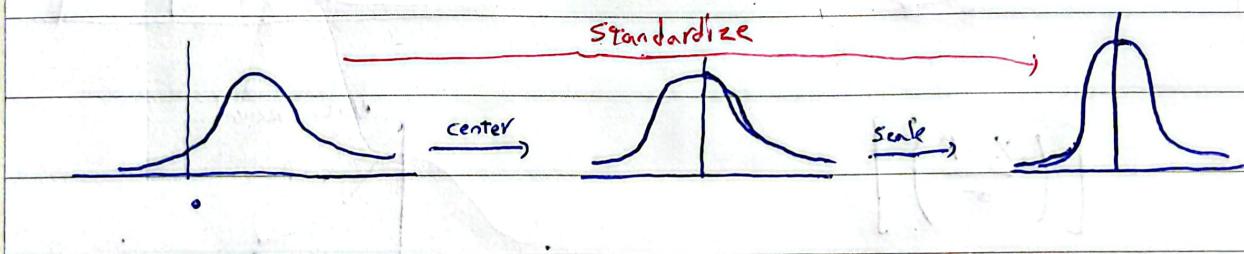
**Sum of Gaussians:** If  $X \sim N(\mu_x, \sigma_x^2)$  and  $Y \sim N(\mu_y, \sigma_y^2)$  are independent, and  $W = aX + bY \Rightarrow W \sim N(a\mu_x + b\mu_y, a^2\sigma_x^2 + b^2\sigma_y^2)$



$$\text{Var}(cX) = c^2 \text{Var}(X) \quad \text{std}(cX) = c \text{std}(X)$$

**Standardize a distribution:** It is good to have  $\mu=0$  and  $\text{std}=1$ .

$$X \xrightarrow{\text{center}} X - \mu \xrightarrow{\text{scale}} \frac{X - \mu}{\sigma}$$



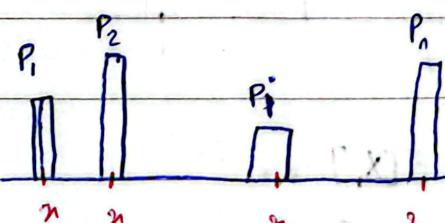
**benefits:** 1- Comparability between different datasets 2- Simplification of statistical analysis  
3- Improved performance of ML models (improve the convergence rate of optimization algorithms and prevent some features from dominating others, leading to improved model performance)

**Moment of distribution**

$$E[X] = P_1 x_1 + \dots + P_n x_n$$

$$E[X^2] = P_1 x_1^2 + \dots + P_n x_n^2$$

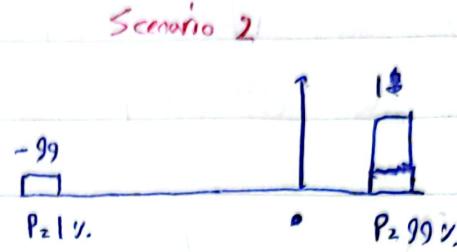
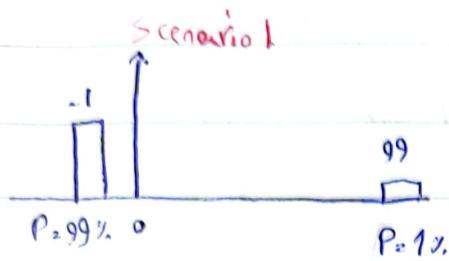
$$E[X^k] = P_1 x_1^k + \dots + P_n x_n^k$$



$X$  - Random var.

Subject: \_\_\_\_\_

Date: \_\_\_\_\_

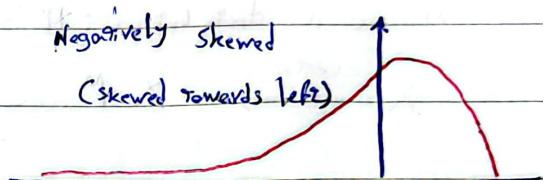
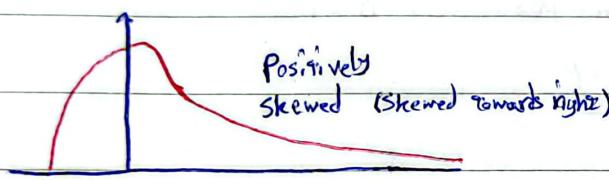


$$E[X_1] = -1 \times 0.99 + 99 \times 0.01 = 0 = E[X_2]$$

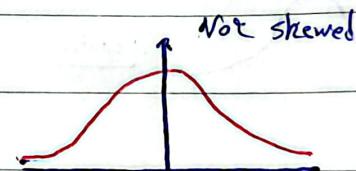
$$E[X_1^2] = (-1)^2 \times 0.99 + (99)^2 \times 0.01 = 99 = E[X_2^2]$$

$$E[X_1^3] = (-1)^3 \times 0.99 + (99)^3 \times 0.01 = 9702 \neq E[X_2^3] = (1)^3 \times 0.99 + (-99)^3 \times 0.01 = -9702$$

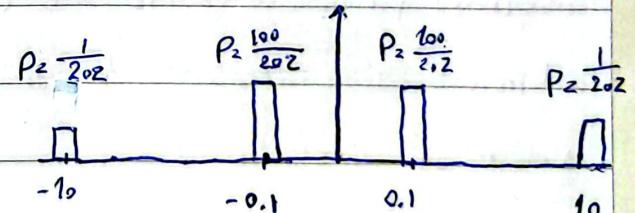
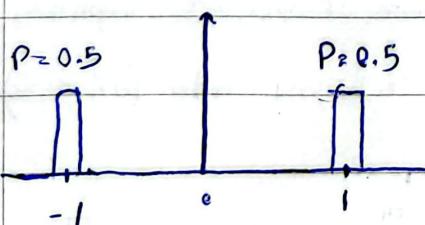
\* The cube of the var detects if you have numbers that are skewed towards the right or skewed towards the left



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



Scenario 1



$$E[X_1] = 0$$

$$E[X_1^2] = 1$$

$$E[X_1^3] = 0 \quad (\text{skew}(X_1) = 0)$$

$$E[X_1^4] = 1$$

2

2

2

≠

$$E[X_2] = 0$$

$$E[X_2^2] = 1$$

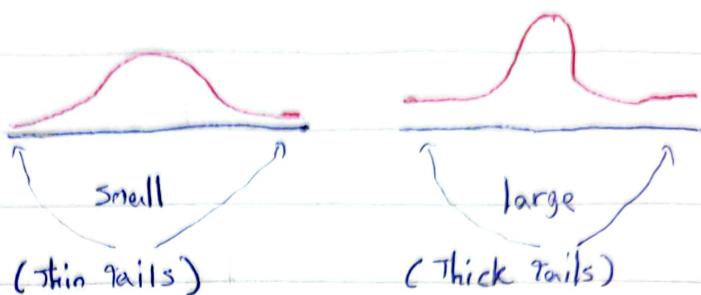
$$E[X_2^3] = 0 \quad (\text{skew}(X_2) = 0)$$

$$E[X_2^4] = 99.01$$

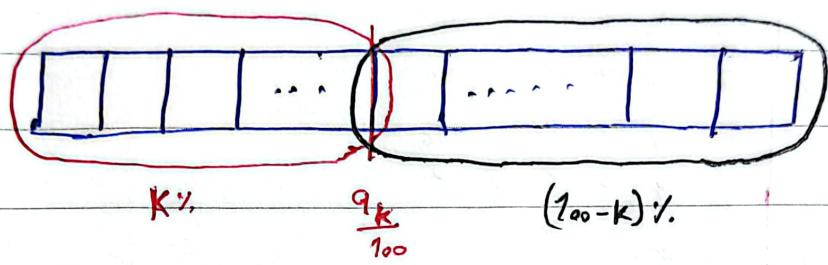
Date:

Subject:

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



**Quantiles:** The  $k\%$  quantile ( $q_{k/100}$ ) is the value that leaves  $k\%$  of your data to the left and  $(100-k)\%$  of your data to the right.



→ some common quantiles: 1. 25% quantile (first ~~quartile~~ - Q1)

2. 50% quantile (median - Q2) . 3. 75% quantile (third quartile - Q3)

\*  $k\%$  quantile ( $q_{k/100}$ ) is the value such that  $P(X \leq q_{k/100}) = \frac{k}{100}$

\* Interquartile range (IQR) =  $Q_3 - Q_1$

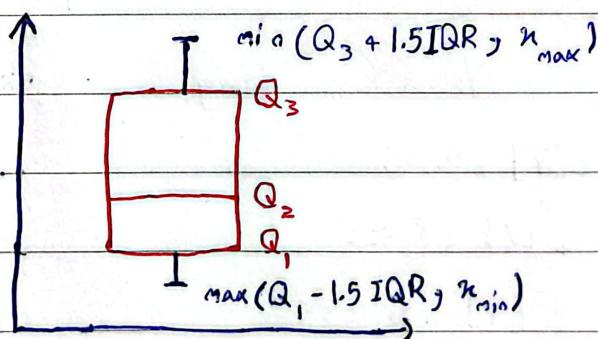
### Box-plots:

Benefits:

- Data is skewed?

- No outlier data? (data were cut at max and min value)

- Analyze dispersion



### Joint Distribution

$$P(x,y) = P(X=x, Y=y)$$

$$P_{XY}(8,48) = P(X=8, Y=48) = 0$$

$$P_{XY}(7,46) = \frac{2}{10}$$

Age  
(X)

|    | 45   | 46   | 47   | 48 | 49   | 50   |
|----|------|------|------|----|------|------|
| 7  | 1/10 | 2/10 | 0    | 0  | 0    | 0    |
| 8  | 0    | 0    | 2/10 | 0  | 0    | 0    |
| 9  | 0    | 0    | 0    | 0  | 3/10 | 1/10 |
| 10 | 0    | 0    | 0    | 0  | 0    | 1/10 |

Subject: \_\_\_\_\_

Date: \_\_\_\_\_

 $X$  = the number rolled on the 1st dice $Y$  = sum of two dice 2nd = $\rightarrow X$  and  $Y$  are independent  $\rightarrow P_{XY}(x,y) = P(x), P(y)$  $X$  = the number rolled on the 1st dice $Y$  = sum of two dice

|    |                |                |                |                |                |                |
|----|----------------|----------------|----------------|----------------|----------------|----------------|
| 12 | 0              | 0              | 0              | 0              | 0              | $\frac{1}{36}$ |
| 11 | 0              | 0              | 0              | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ |
| 10 | 0              | 0              | 0              | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ |
| 9  | 0              | 0              | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ |
| 8  | 0              | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ |
| 7  | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ |
| 6  | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | 0              |
| 5  | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | 0              |
| 4  | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | 0              |
| 3  | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | 0              |
| 2  | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | $\frac{1}{36}$ | 0              |
| 1  | 0              | 0              | 0              | 0              | 0              | 0              |
|    | 1              | 2              | 3              | 4              | 5              | 6              |

$$P(3,7) = \frac{1}{36}$$

$$P(1,1) = 0$$

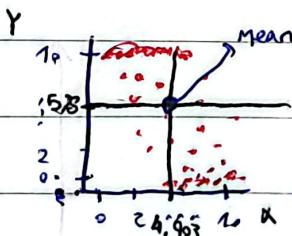
### Joint Continuous Distribution:

 $X$ : waiting time $Y$ : Satisfaction rating

$$E[X] = 4.903 \quad E[Y] = 5.28$$

$$\text{Var}(X) = E[X^2] - E^2[X] = 8.526$$

$$\text{Var}(Y) = E[Y^2] - E^2[Y] = 10.163$$



### Marginal Distribution: reduce the dimensions that we don't need anymore

$$P(Y_j) = \sum_i P_{XY}(x_i, y_j)$$

$$P(X_i) = \sum_j P_{XY}(x_i, y_j)$$

Age(X)

|        |    | Height (Y)     |                |                |    |                |                |
|--------|----|----------------|----------------|----------------|----|----------------|----------------|
|        |    | 45             | 46             | 47             | 48 | 49             | 50             |
| Age(X) | 7  | $\frac{1}{10}$ | $\frac{2}{10}$ | 0              | 0  | 0              | $\frac{3}{10}$ |
|        | 8  | 0              | 0              | $\frac{2}{10}$ | 0  | 0              | $\frac{2}{10}$ |
| Age(X) | 9  | 0              | 0              | 0              | 0  | $\frac{3}{10}$ | $\frac{1}{10}$ |
|        | 10 | 0              | 0              | 0              | 0  | 0              | $\frac{1}{10}$ |

Discrete Conditional Distribution

$$f_{Y|X=n}(y) = \frac{P_{XY}(n, y)}{P_X(n)}$$

e.g.:  $P_{Y|X=9} \text{ (49)} = \frac{3}{4}$

Continuous Conditional Distribution

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

Covariance: Find the relation between two variables

$$\text{Cov}(X, Y) = \frac{\sum (x_i - \mu_x)(y_i - \mu_y)}{n}$$

$\text{Cov}(X, Y) > 0 \rightarrow$  Variables grow together

$\text{Cov}(X, Y) \approx 0 \rightarrow$  Variables are independent from each other

$\text{Cov}(X, Y) < 0 \rightarrow$  when one of them increases, the other one decreases

\* If the probabilities are not equal  $\rightarrow \text{Cov}(X, Y) = \sum P_{XY}(n_i, y_i)(n_i - \mu_x)(y_i - \mu_y)$

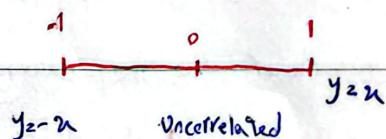
\*  $\text{Cov}(X, Y) = E[XY] - E[X]E[Y]$

Covariance Matrix  $\Sigma$ :

\*  $\text{Cov}(a, b) = \text{Cov}(b, a)$

|   | X                  | Y                  | Z                  |
|---|--------------------|--------------------|--------------------|
| X | $\text{Var}(X)$    | $\text{Cov}(X, Y)$ | $\text{Cov}(X, Z)$ |
| Y | $\text{Cov}(X, Y)$ | $\text{Var}(Y)$    | $\text{Cov}(Y, Z)$ |
| Z | $\text{Cov}(X, Z)$ | $\text{Cov}(Y, Z)$ | $\text{Var}(Z)$    |

Correlation Coefficient:  $\frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)} \sqrt{\text{Var}(Y)}}$



Multivariate Gaussian Distribution

$$f_X(n) = \frac{1}{\sqrt{2\pi}^6} e^{-\frac{1}{2} \frac{(n-\mu)^T}{\sigma^2}}$$

$$f_X(n_1, n_2, \dots, n_n) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} e^{-\frac{1}{2} (n-\mu)^T \Sigma^{-1} (n-\mu)} \quad n = [n_1, n_2, \dots, n_n]^T$$

**Population ( $N$ ):** The entire group of individuals or elements you want to study which share a common behaviour

**Sample ( $n$ ):** Subset of the population you use to draw conclusions about the population as a whole

\* Every dataset you work with in machine learning is a sample

**Proportion:**  $P = \frac{\text{number of items with a given characteristic } (x)}{\text{Population } (n)}$

\* Sample proportion  $\rightarrow \hat{P}$  or  $P'$

$$\text{Population Variance: } \sigma^2 = \frac{1}{N} \sum (n - \bar{M})^2$$

↑ Population mean  
Var(X)      ↑ Population size

$$\text{Sample Variance: } S^2 = \frac{1}{n-1} \sum (n - \bar{x})^2$$

↑ sample mean  
↑ sample size

**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.

(Under certain conditions): 1-Sample is randomly drawn. 2-Sample size must be sufficiently large 3-Independent observations

**Central limit theorem:** As  $n$  (size of sample) increases, the probability distribution becomes closer to a Gaussian Distribution with  $\bar{x} \approx M = np$  and  $\sigma^2 \approx np(1-p)$

Date: \_\_\_\_\_

Subject: \_\_\_\_\_

(CTL - Example 2: Continuous Random Var (Independent))

$$\bar{Y}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\rightarrow \mu_{\bar{Y}_n} = \mu \rightarrow \text{Population}$$

$$\sigma_{\bar{Y}_n}^2 = \frac{\sigma^2}{n} \rightarrow \text{Population Variance}$$

Maximum likelihood estimation: Find model that most likely produced the data or

Maximize  $P(\text{Data} | \text{model})$ Bernoulli Example:  $n$  coins  $X = (X_1, \dots, X_n)$   
 $X_i \stackrel{\text{iid}}{\sim} \text{Bernoulli}(p)$ 

Likelihood  $L(p; n) = P(X=n) = \prod_{i=1}^n P_{X_i}(n_i) = \prod_{i=1}^n p^{n_i} (1-p)^{1-n_i}$  taking log-likelihood derivative  
 and find  $\hat{p} \rightarrow \hat{p} = \frac{\sum_{i=1}^n n_i}{n} = \bar{n}$

Gaussian Example: The best distribution is the one where the mean-variance of the dist is the mean-variance of the sample

Linear Regression Example: Find a line that minimizes the sum of squared distance from the point (covered in Calculus)

|                | Loss     | Equation  | Penalty          | New loss |
|----------------|----------|---|------------------|----------|
| Regularization | $\eta_0$ | $y = 4x + 3$                                    | $L_2 = 4^2 + 16$ | 26       |
| Data           | 2        | $y = 2x^2 - 4x + 5$                             | $L_2 = 20$       | 22       |
| Model 3        | 0.1      | $y = 4x^{10} - 9x^8 - 2x^6 + 3x^5 - 6x^4 - 10x$ | $L_2 = 246$      | 246.1    |

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

Log-loss:  $ll$

$L_2$  Regularization Error:  $a_n^2 + a_{n-1}^2 + \dots + a_1^2$

Regularization parameter:  $\lambda$

Regularized error:  $ll + \lambda(a_n^2 + \dots + a_1^2)$

Maximum Likelihood with Bayes

$$P(\text{Data} | \text{Model}) \cdot P(\text{Model})$$

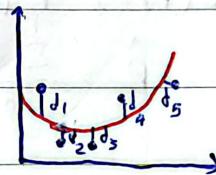
Take logarithms

Polynomial regression with regularization

Square-loss + Regularization term

$$P(\text{Model 1}) \propto \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} a_1^2}$$

$a_1 \sim N(0, 1)$  \*



$$P(\text{Model 2}) \propto \prod_{i=1}^2 \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} a_i^2}$$

$$a_1^2 + a_2^2 + b$$

$$P(\text{Model 3}) \propto \prod_{i=1}^5 \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2} d_i^2}$$

$$a_1 x^1 + \dots + a_{10} x^{10} + b$$

Maximize:  $\log(P(\text{Data} | \text{Model})) + \log(P(\text{Model})) \rightarrow -\frac{1}{2} [d_1^2 + \dots + d_5^2 + a_1^2 + a_2^2]$

Minimize:  $d_1^2 + \dots + d_5^2 + a_1^2 + a_2^2$

## Frequentist Vs Bayesian statistics :

**Frequentists:** 1- Probabilities represent long term frequency of events    2- Concept of likelihood  
 3- Goal: Find the model that most likely generated the observed data

**Bayesians:** 1- Probabilities represent the degree of belief (or certainty)  
 2- concept of prior    3- Goal: Update prior belief based on observations

Maximum A posteriori (MAP-Bayesian):

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

Posterior: A will happen after considering B evidence

Prior: A will happen before considering evidence B

Likelihood of evidence B appearing, given A happened

Probability of evidence B in any circumstances

$P(B|A)P(A) \rightarrow P(B|A')P(A')$

choose the one with highest probability → Make the updated belief Posterior

A: an event you are trying to predict

B: another event, or evidence, that helps refine your prediction

E.g.: We have two types of coin and a mystery coin that is either fair or biased

| Fair              | Biased            | Event to Predict   | Evidence   |
|-------------------|-------------------|--|--|
| $P(H) = 0.5$      | $P(H) = 0.8$      | $A \rightarrow Y$ takes some value   | $B \rightarrow X$ take some value                                  |
|                   |                   | $Y$ : odds of H for your coin.   | $X$ : result of coin flip  |
| $P(Y=0.5) = 0.75$ | $P(Y=0.8) = 0.25$ | $Y = \begin{cases} 0.5 & \text{if coin is fair} \\ 0.8 & \text{if coin is biased} \end{cases}$ | $X = \begin{cases} 0 & \text{if T} \\ 1 & \text{if H} \end{cases}$ |

Priors:  $P(Y=0.5) = 0.75$      $P(Y=0.8) = 0.25$

$$P_{Y|X=1}(0.5) = \frac{P_x(Y=0.5|1) P_y(0.5)}{P_x(1)} = \frac{0.5 \cdot 0.75}{0.5 \cdot 0.75 + 0.8 \cdot 0.25} \rightarrow P(Y=0.5|X=1) = 0.652$$

Posterior

New Prior

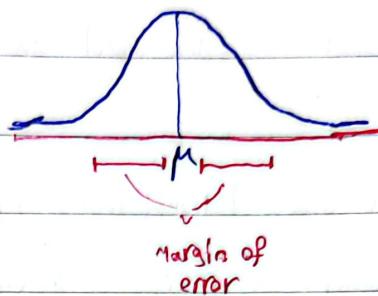
$$P_{Y=0.8|X=1} = 0.348$$

**Confidence Interval:** Interval of values which contains the population parameter (like  $\mu$ ) with some degree of certainty.

$$\text{Unknown} \leftarrow \rightarrow \text{known}$$

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

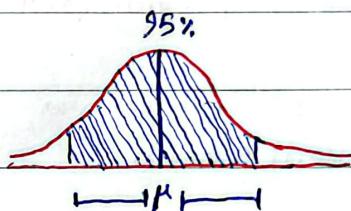
$\bar{X} \sim N(\mu, \sigma^2)$ : Describe the prob. of selecting different  
(of size = 1)  
sample means



\* Confidence level: Prob. that your sample mean is within the margin of error

\* Significance level ( $\alpha$ ):  $1 - \text{confidence level} \xrightarrow{\text{so}}$  confidence level  $= 1 - \alpha$

e.g.:  $\alpha = 0.05 \rightarrow \text{confidence level} = 0.95 \text{ or } 95\%$



→ 95% of all randomly generated sample means would fall within the shaded portion of the distribution curve

confidence interval:  $\bar{x} \pm \text{margin of error}$   
↳ sample mean

\* When confidence level is  $\alpha\%$  →  $\alpha\%$  of confidence intervals would contain  $\mu$

→ What if sample size  $= k$  (instead of 1)  $\rightarrow \frac{\mu}{\bar{x}} = \mu \quad \frac{\sigma}{\sqrt{k}} = \frac{\sigma}{\sqrt{k}}$

\* Increasing sample size

→ Decreasing margin of error

→ Increasing accuracy ( $\bar{x}$  is closer to  $\mu$ )

→ as Margin of error decreases  $\rightarrow$  Confidence interval shrinks

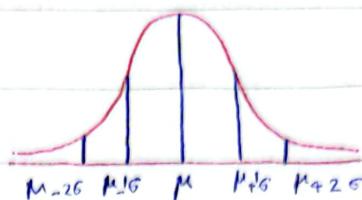
\* Decreasing confidence level  $\rightarrow$  Decreasing margin of error

\* Ideally you want high confidence level and narrow interval  $\rightarrow$  ~~more~~ more data would help!

Margin of error

$$(\mu - 1\sigma, \mu + 1\sigma) \rightarrow 68\%$$

$$(\mu - 2\sigma, \mu + 2\sigma) \rightarrow 95\%$$

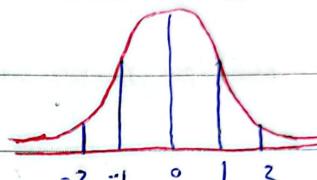


O: Z-score

\* Z-distribution (standard normal distribution):  $\frac{\bar{X} - \mu}{\sigma} = Z \sim N(0, 1^2)$

$$(-1, 1) \rightarrow 68\%$$

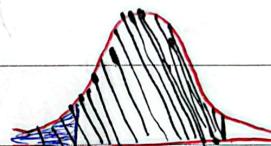
$$(-2, 2) \rightarrow 95\%$$



\*  $(-\underline{1.96}, \underline{1.96})$  is exact point for 95%.

critical values  $Z_{\alpha/2}$  has  $\alpha\%$  to its left

$$Z_{0.975} = 97.5\% \text{ to its left}$$



\* if significant level  $= \alpha$   $\rightarrow Z_{1-\alpha/2} - Z_{\alpha/2}$ ; between these two values, lies  $(1-\alpha)\%$  of the distribution

But if we had non-standardized normal dist, we multiply these critical values by standard deviation  $\rightarrow$  e.g.: 95%:  $(\mu - 1.96\sigma, \mu + 1.96\sigma)$

$$\text{Since } \sigma_{\bar{x}} = \frac{\sigma}{\sqrt{n}} \rightarrow 95\%: (\mu - 1.96 \frac{\sigma}{\sqrt{n}}, \mu + 1.96 \frac{\sigma}{\sqrt{n}})$$

In general: Margin of Error:  $(\mu - Z_{1-\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}, \mu + Z_{1-\alpha/2} \cdot \frac{\sigma}{\sqrt{n}})$

$$\text{So: } \mu - 1.96 \frac{\sigma}{\sqrt{n}} < \bar{x} < \mu + 1.96 \frac{\sigma}{\sqrt{n}} \rightarrow \bar{x} - 1.96 \frac{\sigma}{\sqrt{n}} < \mu < \bar{x} + 1.96 \frac{\sigma}{\sqrt{n}}$$

$$\text{Confidence interval: } \bar{x} \pm Z_{1-\alpha/2} \cdot \frac{\sigma}{\sqrt{n}}$$

### Confidence interval - Calculation Steps:

1- Find sample mean  $\bar{x}$

2- Define a desired confidence level  $(1-\alpha)$

3- Get the critical value  $(Z_{1-\frac{\alpha}{2}})$

4- Find the standard error  $(\frac{\sigma}{\sqrt{n}})$

5- Find the margin of error:  $Z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$

Confidence interval

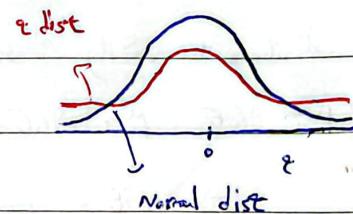
6- Add/subtract the margin of error to the sample mean  $\bar{x} \pm Z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}}$

\* Assumptions: 1-Simple random sample 2-Sample size  $> 30$  or population is approx normal

\*  $n > \left( \frac{Z_{1-\frac{\alpha}{2}} \cdot \sigma}{\text{margin of error}} \right)^2$  : sample size to get the desired MOE  
 (MOE)

Unknown  $\sigma$ :  $\sigma \rightarrow s$  (sample standard deviation)  $\frac{\bar{x} - \mu}{\sigma} \rightarrow \frac{\bar{x} - \mu}{s}$

Student's t dist:



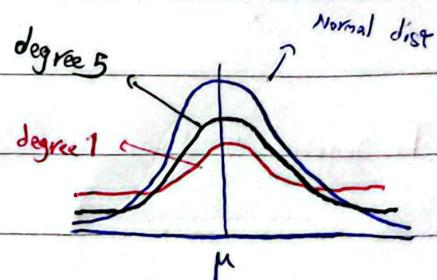
normal dist

t-dist

$$\bar{x} \pm Z_{1-\frac{\alpha}{2}} \cdot \frac{\sigma}{\sqrt{n}} \rightarrow \bar{x} \pm t_{1-\frac{\alpha}{2}} \cdot \frac{s}{\sqrt{n}}$$

sample size

\* There are several t dist! degree of freedom =  $n-1$



\* The larger the number of degrees of freedom  $\rightarrow$  closer to a normal dist

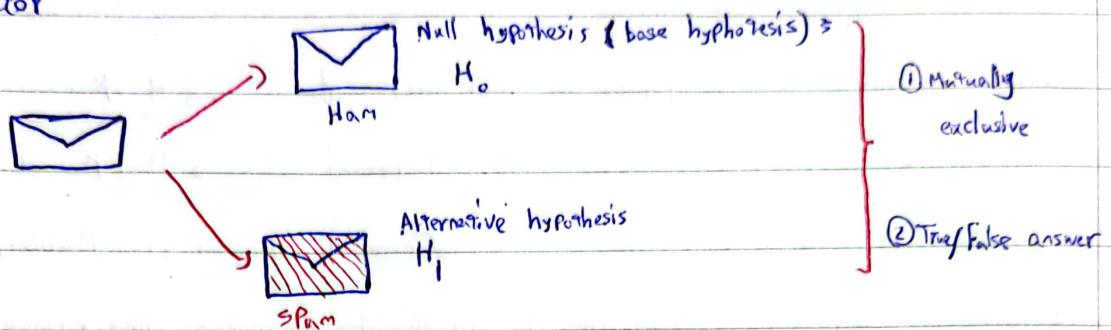
Confidence interval for proportions:

Confidence interval  $\hat{p} \pm \text{margin of error}$

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

**Hypothesis testing:** Is a way to tell if some belief you have of the population is true or not.

e.g. email spam detector



\* If there was alot of evidence that shows that the email is spam  $\rightarrow$  Null hypothesis rejected and alternative is accepted

\* If the evidence gathered is not enough  $\rightarrow$  We are not accepting that the email is ham!

\* we decide between two hypothesis base on: Data - Evidence

### Type I and Type II Errors

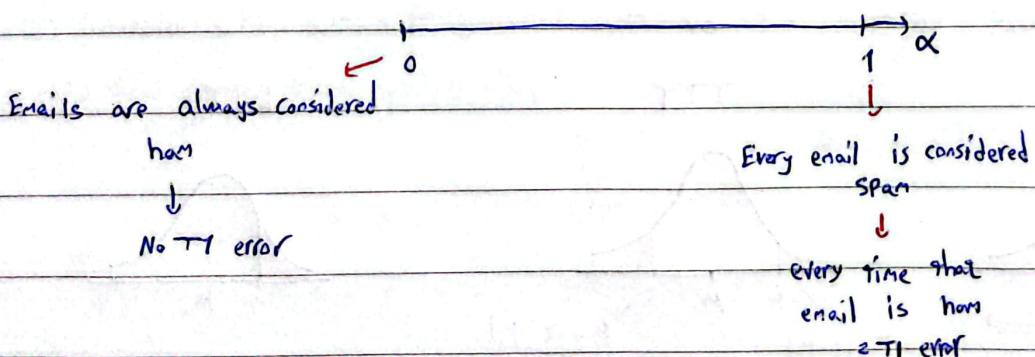
T1  $\rightarrow$  False Positive

T2  $\rightarrow$  False Negative

| Decision           | Reality    |             |
|--------------------|------------|-------------|
|                    | $H_0$ True | $H_0$ False |
| Reject $H_0$       | Type I     | ✓           |
| Don't Reject $H_0$ | ✓          | Type II     |

\* In this case, Type I errors are worse than Type II

**Significance level:** what is the greatest prob of Type I error you are willing to tolerate?



\*  $\alpha = 0.05$  is a good choice

\* ↓ Type I error  $\rightarrow$  ↑ Type II error  
(Too much)

$$\rightarrow \alpha = \max P(T1 \text{ error}) = \max P(\text{Reject } H_0 | H_0)$$

**Data Quality:** Data should be Reliable

→ Representative (of the population)

→ Randomized (No bias)

→ Sample size (consider 30 samples or more)

e.g.: The mean height for 18 yrs in the US in the 70s was 66.7 in, what about now?

3 questions

3 set of hypothesis

- ① Right-Tailed Test: Population mean increased  $\rightarrow H_0: \mu = 66.7$  vs  $H_1: \mu > 66.7$
- ② Left-Tailed Test: - decreased  $\rightarrow H_0: \mu = 66.7$  vs  $H_1: \mu < 66.7$
- ③ Two-Tailed Test: - changed  $\rightarrow H_0: \mu = 66.7$  vs  $H_1: \mu \neq 66.7$

\* hypothesis are always formulated in term of the population (like population mean)

\* decisions = based on the observations Here: ~~at~~ sample mean  $\bar{X} = \frac{1}{n} \sum_{i=1}^{10} X_i$

$\bar{X}$   $\rightarrow$  test statistic: a function of random samples  $T(\bar{X}) = \bar{X} = (X_1, \dots, X_n)$

T1: Determine  $\mu > 66.7$  when  $\mu = 66.7$

① RTT:  $T_1: \mu > 66.7$  when  $\mu = 66.7 \Rightarrow \mu > 66.7$

T1: Determine  $\mu < 66.7$  when  $\mu = 66.7$

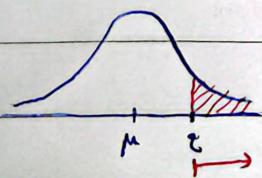
② LTT:  $T_2: \mu < 66.7$  when  $\mu = 66.7 \Rightarrow \mu < 66.7$

T1: Determine  $\mu \neq 66.7$  when  $\mu = 66.7$

③ TTT:  $T_3: \mu \neq 66.7$  when  $\mu = 66.7 \Rightarrow \mu \neq 66.7$

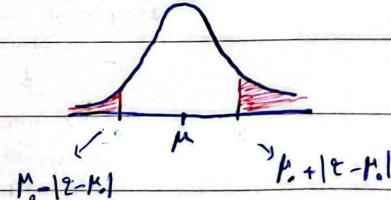
p-values: A p-value is the probability, assuming  $H_0$  is true, that the test statistic takes on a value as extreme as or more extreme than the value observed ( $t$ )

RTT



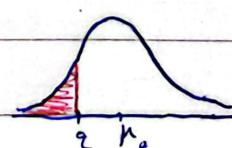
$$P(T(X) > t | H_0)$$

TTT



$$P(|T(X) - \mu_0| > |t - \mu_0| | H_0)$$

LTT



$$P(T(X) < t | H_0)$$

→ if we have  $\bar{X} = 68.492$  and the goal is T1 error prob (P-value)  $< \alpha$

$$\begin{matrix} n = 3 \\ n = 10 \end{matrix}$$

(observed value)

→ If P-value  $< \alpha \rightarrow$  reject  $H_0$  (and accept  $H_1$ )

→ If P-value  $> \alpha \rightarrow$  don't reject  $H_0$ .

RTT:

$$P(\bar{X} > 68.442 \mid \mu = 66.7) = 0.0332 < \alpha = 0.05 \rightarrow \text{Reject } H_0 \text{ (with 5% significance level)}$$

TTT:

$$P(|\bar{X} - 66.7| > |68.442 - 66.7| \mid \mu = 66.7) = 0.0663 > \alpha \rightarrow \text{Don't reject } H_0.$$

LTT:

$$P(\bar{X} < 64.252 \mid \mu = 66.7) = 0.0049 < \alpha \rightarrow \text{Reject } H_0.$$

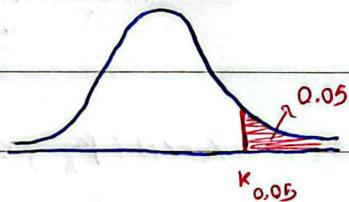
\* If  $H_0$  is true:  $\bar{X} \sim N(\mu_0, \frac{\sigma^2}{n})$  or  $Z = \frac{\bar{X} - \mu_0}{\sigma/\sqrt{n}} \sim N(0, 1)$  (z-statistic)

Critical value: sample that has p-value  $\leq \alpha$  (least extreme sample you could get to reject  $H_0$ )

How to compute:  $\alpha = 0.05$  (Height example)  $0.05 = P(\bar{X} > K_{0.05} \mid \mu = 66.7)$

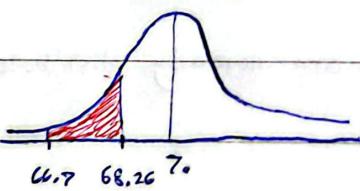
$$K_{0.05} = 68.26$$

Decision rule: Reject  $H_0$  if  $\bar{X} > 68.26$



\* What is the Type II error prob if the true value is  $\mu = 70$ ? If  $\mu = 70 \rightarrow \bar{X} \sim N(70, \frac{\sigma^2}{n})$

$$\underbrace{P(\bar{X} < 68.26 \mid \mu = 70)}_{\beta} = 0.0333$$



Type II error:  $P(\text{Don't reject } H_0 \mid \mu = 70)$

Power of the test:

Power of the test:  $P(\text{Reject } H_0 \mid \mu = 70)$

} complementary  $\rightarrow \text{Power} = 1 - \beta$

Steps for Performing Hypothesis testing:

1. State your hypotheses

a) Null hypothesis: The baseline  $\rightarrow H_0: \mu = 66.7$

b) Alternative hypothesis: the statement you want to prove  $\rightarrow H_1: \mu > 66.7$

## 2- Design your test

a) Decide the test statistic to work with  $\rightarrow \bar{X}$

b) Decide the significance level  $\rightarrow \alpha = 0.05$

3- Compute the observed statistic (based on your sample)  $\bar{x} = 68.492$

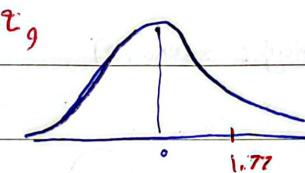
## 4- Reach a conclusion:

a) if p-value  $< \alpha \rightarrow$  reject  $H_0$

t-dist

If  $\sigma$  is unknown: If  $H_0$  is true  $\rightarrow \bar{T} = \frac{\bar{X} - 66.7}{S/\sqrt{n}} \sim t_9$  (t-dist with  $v=9$  (degree of freedom))

repeat above process with  $t_9$ :



Independent Two-Sample t-Test:  $X \sim N(\mu_A, \sigma_A^2)$   $Y \sim N(\mu_B, \sigma_B^2)$

\* Assumptions: 1- All elements in the sample from the two group are different

2- Each element in both samples are independent

3- Populations are normally distributed

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\bar{X} - \bar{Y} \sim N(\mu_A - \mu_B, \frac{\sigma_A^2}{n} + \frac{\sigma_B^2}{m})$$

$$\bar{Y} = \frac{1}{m} \sum_{i=1}^m Y_i$$

if  $\sigma$  is known:

$$\frac{(\bar{X} - \bar{Y}) - (\mu_A - \mu_B)}{\sqrt{\frac{\sigma_A^2}{n} + \frac{\sigma_B^2}{m}}} \sim N(0, 1)$$

if  $\sigma$  is unknown:

$$T = \frac{(\bar{X} - \bar{Y}) - (\mu_A - \mu_B)}{\sqrt{\frac{s_A^2}{n} + \frac{s_B^2}{m}}} \sim t_{\nu}$$

**Paired T-Test:** Difference between samples:  $\bar{D} = \frac{D_1 + D_2 + \dots + D_n}{n}$   $D_i = X_i - Y_i$   
 $D_i \sim N(\mu_D, \sigma_D^2)$

$$\frac{\bar{D} - \mu_D}{\sigma_D / \sqrt{n}} \sim N(0, 1^2) \quad \text{--- } \sigma_D \text{ is unknown} \rightarrow \sigma_D = s_D \rightarrow T = \frac{\bar{D} - \mu}{s_D / \sqrt{n}} \sim t_{n-1}$$

**Z-test for Proportions:** •  $p$  is the population proportion of individuals in a particular category

- $p_0$  is the population proportion under the null hypothesis
- $n$  is the observed number of individuals in the sample from the specified category
- $n$  is the sample size
- $\hat{p} = \frac{n}{n}$  is the sample proportion

$$\rightarrow Z = \frac{\frac{n}{n} - p_0}{\sqrt{p_0(1-p_0)}} \sim N(0, 1)$$

**Two Sample Test for Proportions:** •  $p_1 - p_2$  is the dif in the population proportion between two groups

- $n_1$  is the observed number of individuals in the sample from the specified category from group 1
- $n_2$  is the observed number of individuals in the sample from the specified category from group 2
- $n_1$  is the sample size for group 1
- $n_2$  is the sample size for group 2

$$Z = \frac{\frac{X}{n_1} - \frac{Y}{n_2} - 0}{\sqrt{\frac{(X+Y)(1-\frac{X+Y}{n_1+n_2})}{n_1 n_2}}} \sim N(0, 1^2)$$