

- Head first statistics Dawn Griffiths
- Introductory statistics Neil A Weiss

Chap 1 → Introductory statistics

- Collecting data → summarize → visualization → analysis & making inferences
- Descriptive statistics v/s Inferential statistics
- Abt summarizing data set & visualization
- Prediction of some data in data set

Ex Flower classification → (Inferential statistics)

- distinguish rose & lotus 50 roses & 50 lotus
- we have petal length & width for each recorded
- we now have a new flower given petal length = a units & petal width = b units, determine species of flower
- We find object with nearest data and this flower belongs to this category.

↓  
optimize

- find centroid of data for features & find nearest category

Mean   Median   Mode

- Just compare it with centroid of diff. category & interpret.

Descriptive Statistics →

- Measures of central tendency → mean, median, mode
- Measure of variance → range, IQR, percentile.
- Data visualization → piechart, Box plot, Bar charts, stem & leaf diagram, surface plots. → use matplotlib library in python

- Observational study  $\rightarrow$  based on observation only
- Experimental study  $\rightarrow$  based on experiment done on some data grp

### Inferential Statistics $\rightarrow$

- regression - relationship b/w dependent & independent variable
- naive Bayes classification

$$y = w_0 + w_1 n_1 + w_2 n_2^2$$

↳ Polynomial regression of order 2

Ex - Need to predict age given value of Blood sugar, weight, height  
 $\therefore$  we need to establish a function for that we make a assumption known as hypothesis

let Age =  $w_0 + w_1 \cdot BS + w_2 \cdot \frac{Weight}{n_1} + w_3 \cdot \frac{height}{n_2}$   
 or  $y$

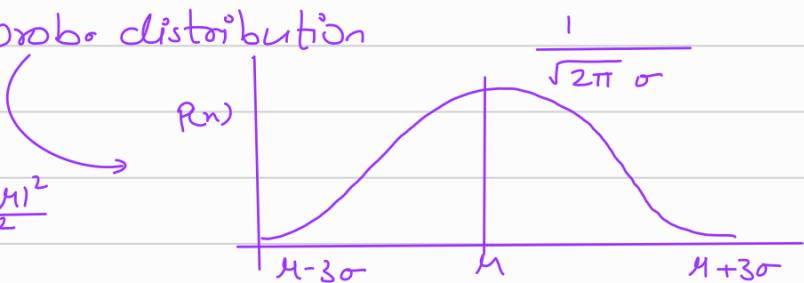
let  $y = w_0 + w_1 n_1 + w_2 n_2 + w_3 n_3 + w_4 \cdot n_1^2 + w_5 n_2^2 + w_6 n_3^2$   
 $+ w_7 n_1 n_2 + w_8 n_2 n_3 + w_9 n_1 n_3$

- optimal order = 1 more than number of terms
- sample should be with no bias & ideal group should be selected from population
- Sampling technique  $\rightarrow$  random sampling, systematic random sampling, cluster based sampling.

Discrete variable  $\rightarrow$  table (Prob. dist.)

- Random variable  $\rightarrow$  prob. distribution

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



 Tally marks :-

Ex 11 22 23 33 44 45 55 11 22 35 6  
 One pass

1      ||||  
2      |||||  
3      ||||  
4      |||  
5      |||  
6      |

- Write a prog to compute mean & standard deviation ( $\sigma$ ) ?
- We just need 1 pass.

$$\sigma^2 = \frac{\sum (n_i - \mu)^2}{n} = \frac{\sum (n_i^2 - 2n_i\mu + \mu^2)}{n}$$

$$= \frac{\sum_{i=1}^N n_i^2}{N} - 2\mu \frac{\sum_{i=1}^N n_i}{N} + \frac{N\mu^2}{N}$$

$$= \frac{\sum_{i=1}^N n_i^2}{N} - 2\mu^2 + \mu^2$$

$$= \frac{\sum_{i=1}^N n_i^2}{N} - \mu^2$$

$$= \frac{\sum_{i=1}^N n_i^2}{N} - \left( \frac{\sum n_i}{N} \right)^2$$

 Computation formula for standard deviation.

- If we use above method upon adding data our previous computation are not lost

- An algorithm is said to be incremental when u can update the model on availability of additional data.

Design experiment →

- Experimental units
- treatment → 8
- factor → Irr. regin, polymer
- level of factor → Irr. regin ⇒ 4, polymer ⇒ 2
- response variable → wt. gain

Er Impact of irrigation regin & use of a particular polymer? for max weight gain of cactus plant

irrigation regin → none low moderate high  
 polymer → yes no

	N	L	M	H
Y	T1	T2	T3	T4
N	T5	T6	T7	T8

0. A ML algorithm named SPM has 2 free parameters named regularization parameter and kernel width parameter,  $\sigma$ . In order to find optimal values of parameter we perform a grid search

$$C \in \{2^{-18}, 2^{-16}, \dots 2^0, \dots 2^{50}\} = 35 \text{ terms}$$

$$\sigma \in \{2^{-18}, 2^{-16}, \dots 2^0, \dots 2^{20}\} = 20 \text{ terms}$$

for which accuracy is maximum

factors → C,  $\sigma$

level of factor → C = 35,  $\sigma$  = 20

Response var → Accuracy

treatment → 35 × 20

- Q. An exp. was conducted to study the impact of folic acid on birth defect, to perform this study a grp of 200 women considered, 100 women were given 10 mg folic acid tablet & the rest of women were given trace element find out the factors, number of treatment experimental variable & response variable!

Exp. unit  $\Rightarrow$  subject (Biology) 200

Res. var  $\Rightarrow$  child birth

factor  $\Rightarrow$  Acid

level  $\Rightarrow$  2

Treatment  $\Rightarrow$  2

Sampling  $\rightarrow$  with replacement  
 $\rightarrow$  Without replacement

How to distribute experimental unit among treatments  $\Rightarrow$

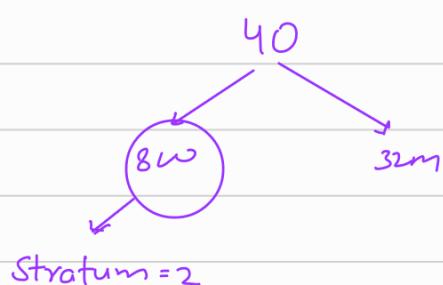
- Complete randomized design
- Randomized block design

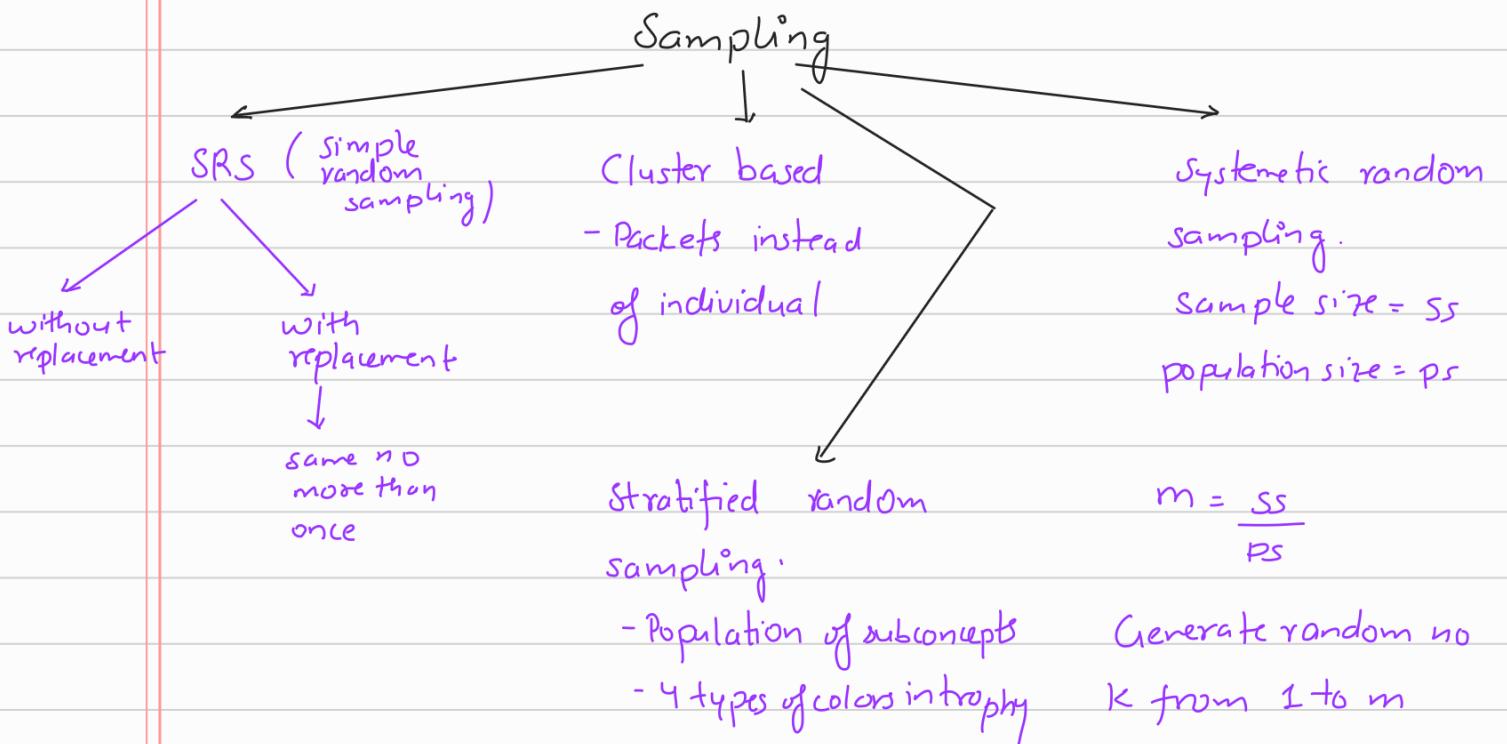
- Q. Driving distance of golf ball - 4 Brands golf balls  
 there are 40 golfers 8 are women & 32 are men

GB1	GB2	GB3	GB4
↓ 10	↓ 10	↓ 10	↓ 10
8M	10M	10M	10M

GB1 driving dist less  
 As sample is not good.

- Sample should be representative of population
- For good same thirshld be  $8M + 2W$





Now my samples are  
 $k, k+m, k+2m, \dots$

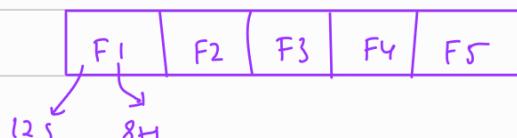
- Cross validation used stratified random sampling
- 



80 training  
20 testing

Then how??

data of 5 equal parts & maintain ratio



Train

Test

T1:	F1	F2 F3 F4 F5
T2:	F2	F1 F3 F4 F5
T3:	F3	F1 F2 F4 F5
T4:	F4	F1 F2 F3 F5
T5:	F5	F1 F2 F3 F4

$m$ -classes

$c_1 \quad n_1$

$c_2 \quad n_2$

$c_3 \quad n_3$

$\vdots \quad \vdots$

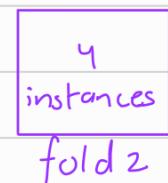
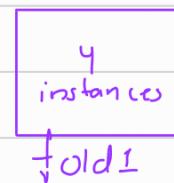
$c_m \quad n_m$

Find unique no. of experimental set up as per using  $k$ -fold cross validation.

$$\frac{n_1 + n_2 + n_3 + n_4}{k} = \frac{y_1 + y_2 + y_3 + y_4}{L} = 8 \text{ instances}$$

2-fold cross validation

↳ divide in 2 equal part



maintaining ratio.

Training

fold 2

fold 1

Test

fold 1

fold 2

- I/P in machine learning is data

- In regression target var. continuous & in classification target is discrete

Training

T  $n$

S  $s'$

T  $t'$

T  $t'$

S  $y'$

Test

P T n

S S  $s'_{1''}$

S S  $s'_{2''}$

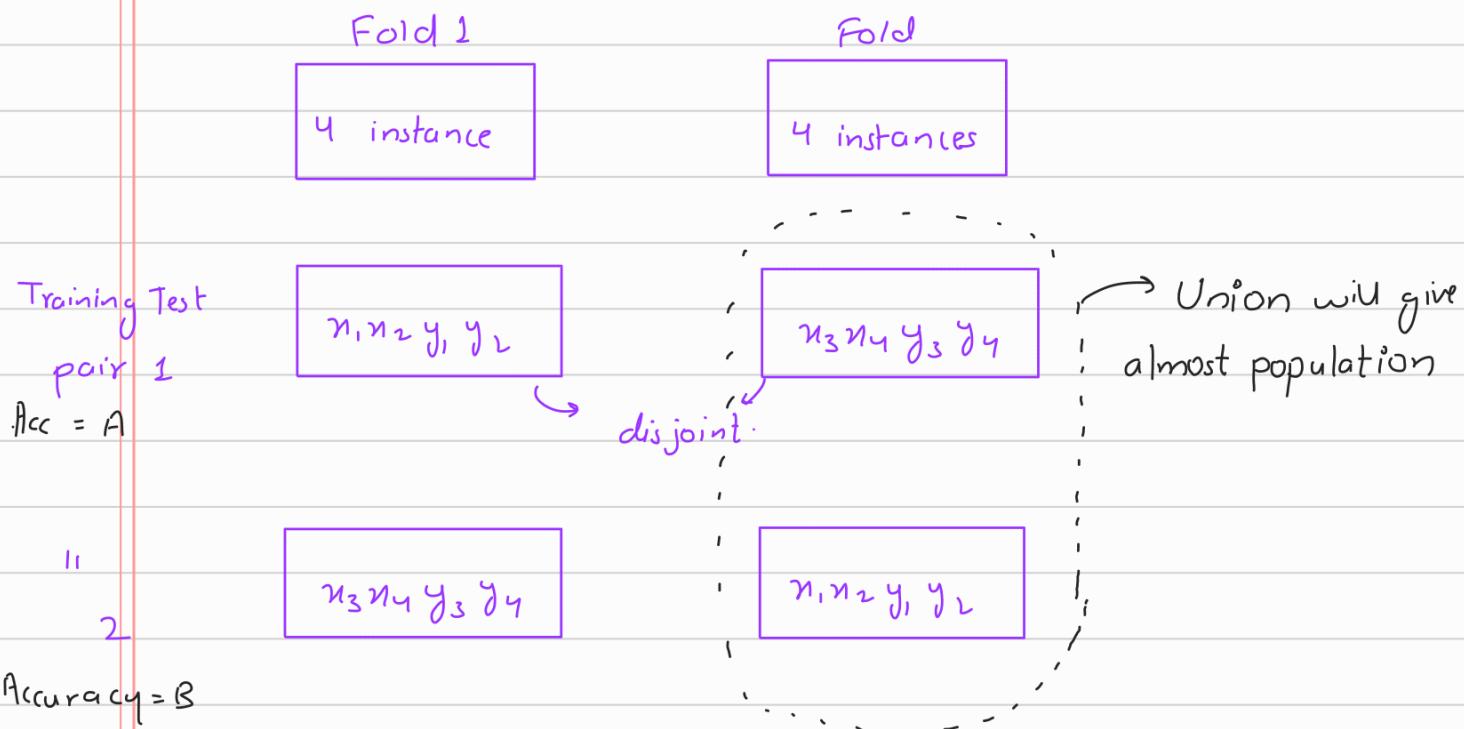
T T  $t'_{3''}$

T T  $t'_{7''}$

$s^* 4.5' \quad t^* 6.5'$

$s^* 5' 1.5'' \quad t^* 6' 5''$

- Here test instance is easy
- To get better do we use k (2) fold cross validation



$$\text{Net Accuracy} = \frac{A+B}{2}$$

$$\text{Unique} = \frac{n_{c_1} \times n_{c_2}}{2}$$

Ass.  $\rightarrow$

$n_1, \dots, n_{10} \Rightarrow \text{Class 1}$

$y_1, \dots, y_{10} \Rightarrow \text{Class 2}$

find out unique experimental setups using 2-folds cross validation  
 prove that no. of unique experimental setups for k-fold cross validation where the  $c_1, c_2, \dots, c_m$  are the  $m$  classes with  $n_1, n_2, \dots, n_m$  instances respectively is

$$\frac{\prod_{i=1}^m n_i}{k} \leq n_i/k$$

Not suitable

Mean

In presence of outliers

(Teacher taught class)

Medians

(Parent with kid in swimming class)

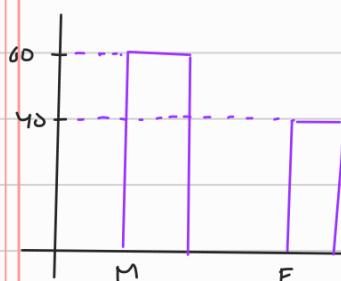
Mode

(5 5 36 36)

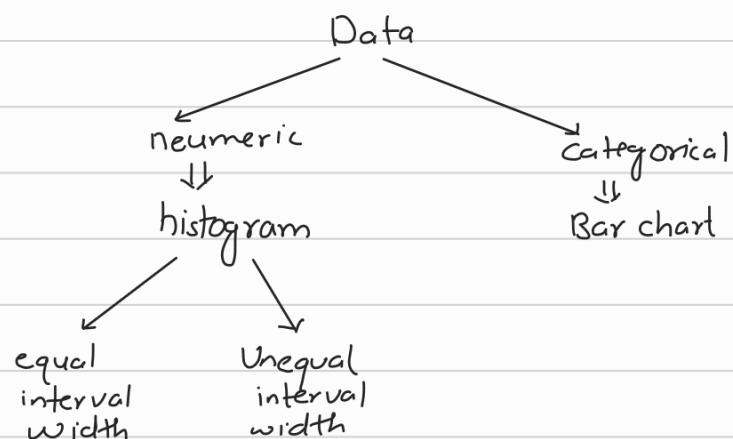
$$\text{mode} = \frac{36+5}{2} = 20.5$$

Distribution →

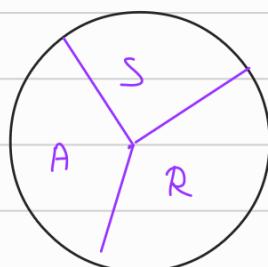
- bar-chart →



Height of bar  $\propto$  Freq.

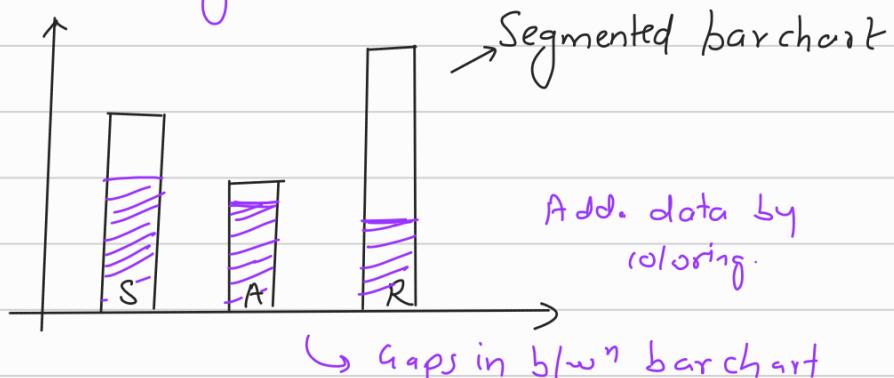


- Pie-chart →



How to plot additional data

- Using barcharts



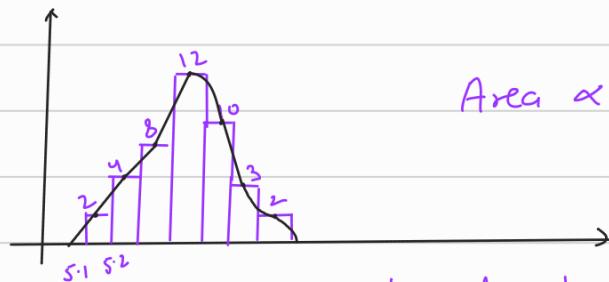
- Histogram → freq

5 → 5.1 2

5.1 → 5.2 4

5.2 → 5.3 8

$6.1 \rightarrow 6.2$       8  
 $6.2 \rightarrow 6.3$       3



Area  $\propto$  Freq.

- No gap b/w<sup>n</sup> bar

- to reduce height we can use fraction
- join centre of histogram

Histogram  $\rightarrow$

3000 Customers

0 $\rightarrow$ 2	1000
2 $\rightarrow$ 4	500
4 $\rightarrow$ 8	500
8 $\rightarrow$ 24	1000
lower limit	
Upper limit	
[8, 24)	

Bar chart



B > A ?? No Bar chart drawback  
(For Nominal or Categorical attribute)

Histogram

Class	freq	Class width	height
0 $\rightarrow$ 2	1000	2	500
2 $\rightarrow$ 4	500	2	250
4 $\rightarrow$ 8	500	4	125
8 $\rightarrow$ 24	1000	16	62.5



Box plot  $\rightarrow$

Central tendency

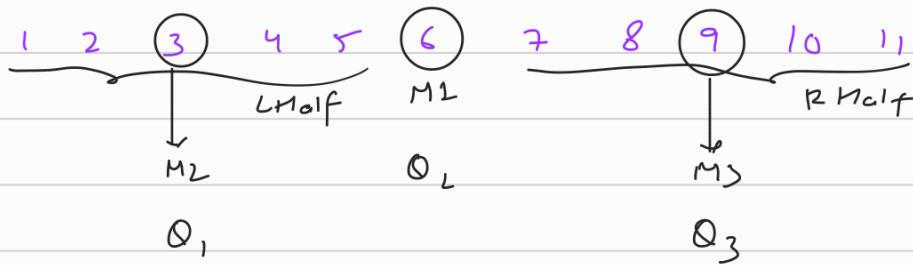
- Mean (fail in outlier)
- Median (fail <sup>kid + parent</sup> major class)
- Mode

Variance

- trimmed range (Range - Outlier)
- range (It contains outlier)
- IQR (for median) - Interquartile range  $\Rightarrow 25\% \quad 50\% \quad 75\%$   
 $Q_1 \quad Q_L \quad Q_3$   
 $IQR = Q_3 - Q_1$
- Variance

- both collectively ct & var. provide much more info about data

IQR →



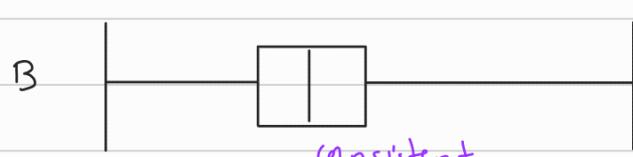
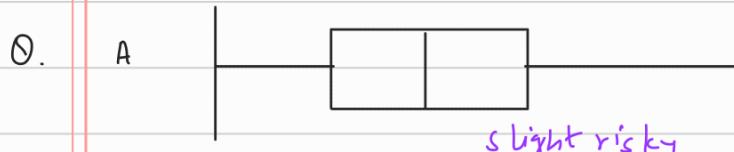
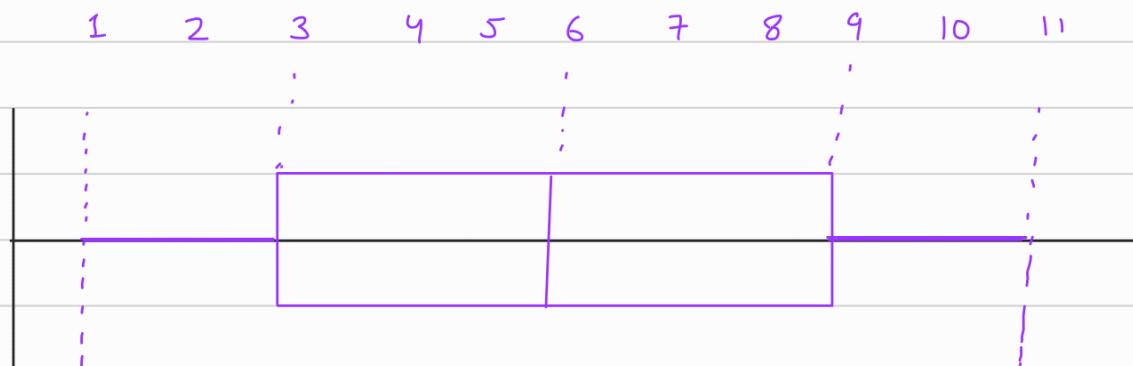
25% data lie below  $Q_1$

lower bound = 1

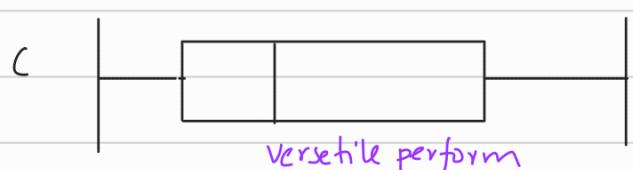
Upper bound = 11

$Q_1, Q_2, Q_3 = 3, 6, 9$

Median = 6



Whom to select ??



Q. Compare the profit of 5 companies for the month Jan - Apr 2023.

- What to use

Piechart

Barchart

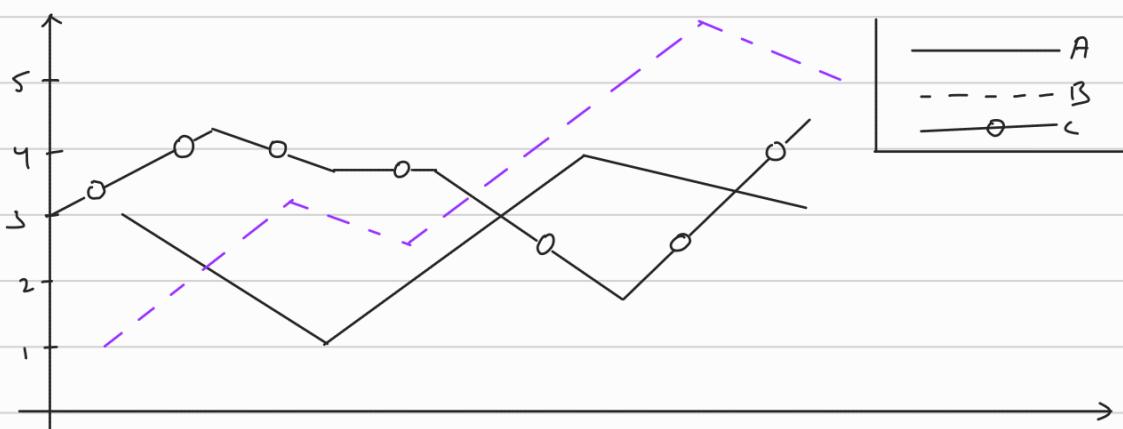
best

line Chart → (diff line similar to below one)

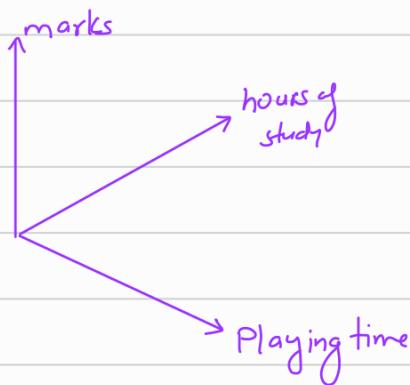
Scatter plots → (diff symbol) for diff people

Just plot point on x-y  
& we create a box on right corner called  
legend declaring symbol symbolization.

Line



Surface plot 3D →



Normal distribution →

- Can be used to summarize numerical data.

Parametric study      (Assume data follow a particular prob. density function)  
Non-parametric

Categorical/Nominal ⇒ prob. distribution table (coin HT table)  
data  
Numeric ⇒ prob. density function

Temp-play →

Yes	No		Hot	Mild	Cool	Legend
Mild	Mild	yes (9)	2/9	4/9	3/9	→ K
Mild	Mild	No (5)	3/5	2/5	9/5	↓ G
Mild	hot					↑ H
Hot						↔ L

summarize

Hot  
cool  
cool

hot  
hot

$$P\left(\frac{\text{Hot}}{\text{Yes}}\right) ??$$

## Numeric

Temp - Play

Yes      No

31      41

32      42

33      43

⋮      44

39      45

Assume data follow normal or gaussian distribution

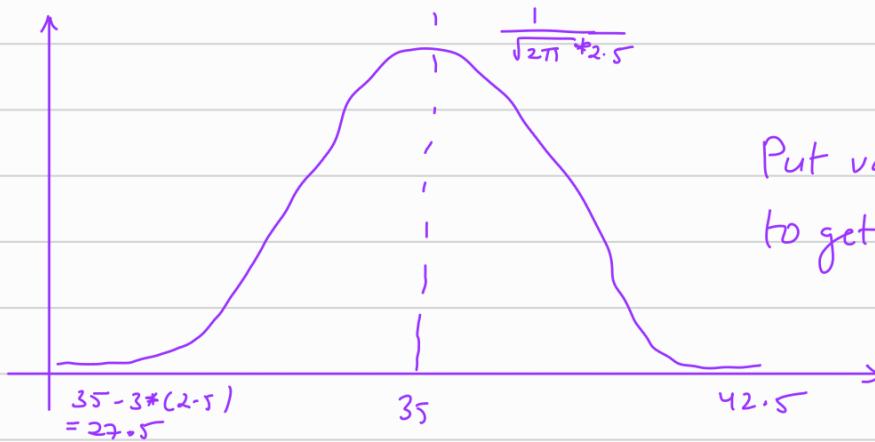
$$f(u) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(u-\mu)^2}{2\sigma^2}}$$

$$\mu, \sigma = ??$$

$$\Rightarrow \text{temp} = 36.5 \quad \text{play} = ??$$

$$\mu = 35, \sigma = \sqrt{\frac{60}{9}} = \frac{7.745}{3} = 2.5$$

$X$  = random var.  
 $n$  = value of random var.



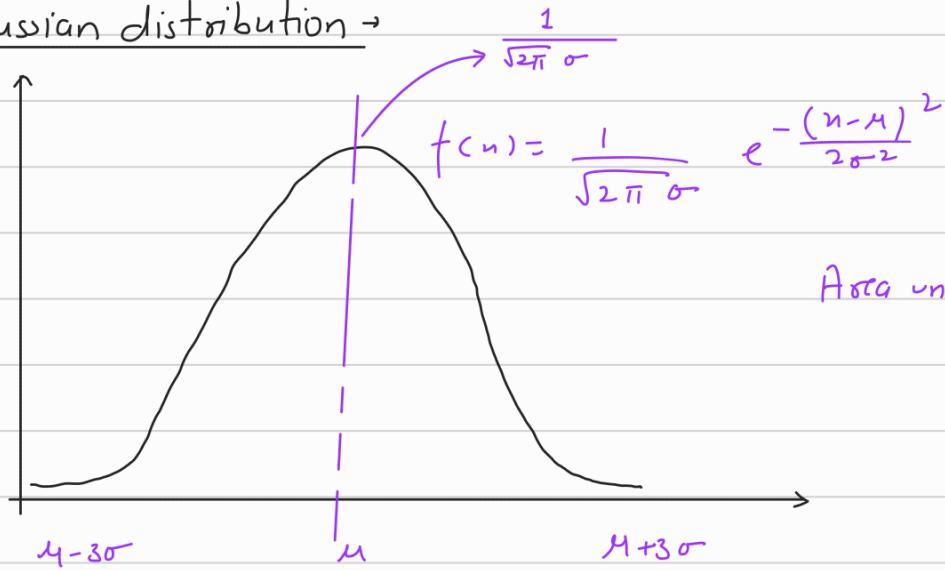
Put value in function  
to get ans.

- Q. Compute the likelihood of play = yes & play = NO when temperature is 40°C assume that the data follows normal/gaussian pdf

$$\rightarrow P(X < 40 \text{ & Play} = \text{Yes}) = \text{Area under curve till } n = 40$$

- To solve Area under curve we convert this distribution into standard normal distribution ( $\mu = 0, \sigma = 1$ )

Gaussian distribution  $\rightarrow$



Area under curve = 1

Empirical rule  $\rightarrow$  (68 - 95 - 99.7 rule)

68% of entire population lie in range of  $\mu - \sigma$  to  $\mu + \sigma$

95% of entire population lie in range of  $\mu - 2\sigma$  to  $\mu + 2\sigma$

99.7% of entire population lie in range of  $\mu - 3\sigma$  to  $\mu + 3\sigma$

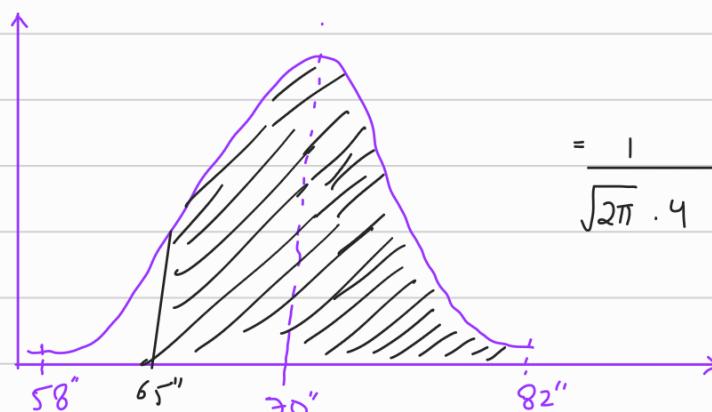
Standard normal distribution  $\rightarrow$  (Z-table)

Mean = 0, Standard deviation = 1

$$f(n) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(n-\mu)^2}{2\sigma^2}} \xrightarrow{\mu=0, \sigma=1} \frac{1}{\sqrt{2\pi}} e^{-\frac{n^2}{2}}$$

eqn of standard normal distribution

- Q. She wants to marry a man who is taller than her. The dist. of the height of men is given by  $N(70'', 16'')$  Given girl height = 65".



integrate in range

or  
transform it to standard normal distribution

$$n = 14$$

$$1 = -4.5$$

to convert  $\sigma$  to 1

$$2 = -3.5$$

divide  $n-1/\sigma$

$$3 = -2.5$$

$$4 = -1.5$$

$$5 = 0.5$$

$$6 = -0.5$$

$$7 = 1.5$$

$$8 = 2.5$$

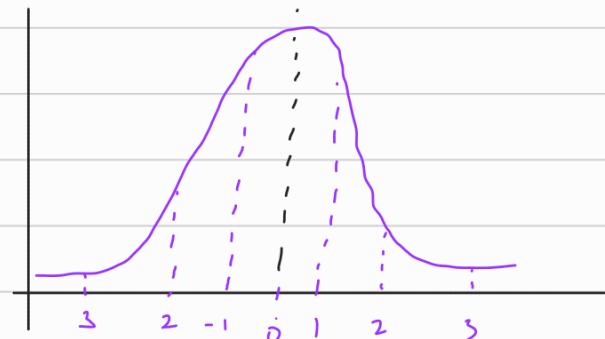
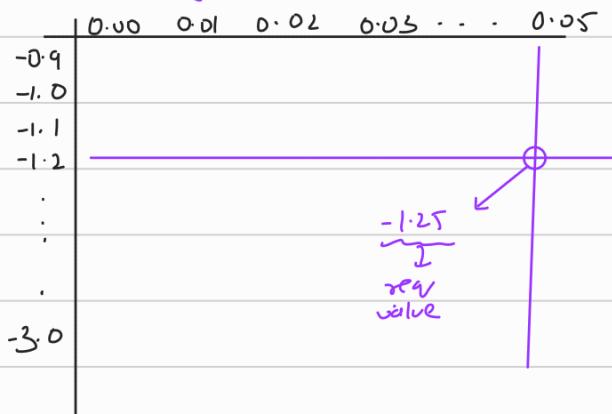
$$9 = 3.5$$

$$10 = 4.5$$

$$\mu = 5.5 \quad \sigma = 0$$

$$z = \frac{65 - 70}{4}$$

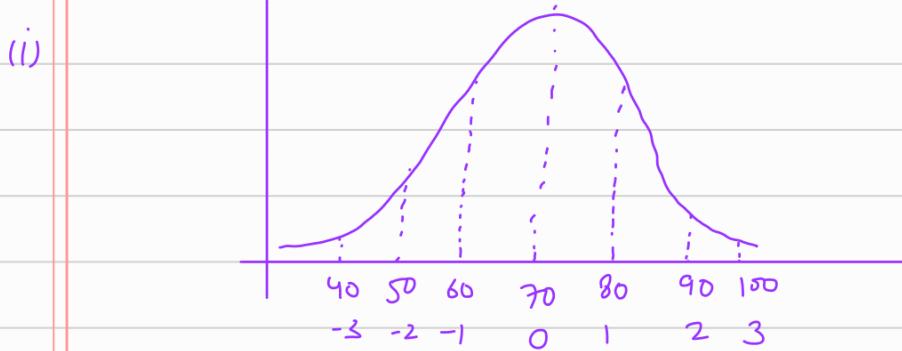
- z-table give area on left hand side of curve



- This will give area upto this portion
- We need area on right  $\therefore 1 - \text{Area of left.}$

- Q. The marks of students of CSE follow gaussian distribution with mean 70 marks & variance of 100 marks find out the following prob.  
 $P(X \geq 80)$ ,  $P(60 \leq X \leq 80)$ ,  $P(X \leq 50)$

-  $\mu = 70, \sigma = 10$



$$P(X \geq 80) = 1 - (\text{z at } 1)$$

$$P(60 \leq X \leq 80) = (\text{z at } 1) - (\text{z at } -1)$$

$$P(Y \leq 50) = (\text{z at } -2)$$

Binomial distribution →

$$P(X = r) = {}^n C_r p^r q^{n-r}$$

$$P(X = 2) = {}^n C_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^2$$

$n=7$

10 in 10

Head 1, 2, 3, 4  
first

Distribution to learn

Geometric

Normal or gaussian

Poisson

Binomial

- O. There are 100 questions with 4 possible ans A,B,C,D in a test all the options are equally likely to be correct , each question carry 1 mark , find out the probability that the no. of ans successfully marks < 70 .



Binomial is over calculating ∴ we will transform it to other distribution

Expectation & variance of random variable:-

E- lottery system → 3-window & D1 symbol → \$, C, L, other

0.1 0.2 0.2 0.5

same for D2 & D3 , say if we get \$ \$ \$ we win 20Rs , \$ \$ C = 15Rs (in any seq.) , CCC = 10Rs , LLL = 5Rs we need to find expected amount we will win , entry fee - 1Rs

- we need to make prob dist. table

X	Prob. ( $X=n$ )
19	$0.1 * 0.1 * 0.1 = 0.001$
14	$0.1 * 0.1 * 0.2 * 3 = 0.006$
9	$0.2 * 0.2 * 0.2 = 0.008$
4	$0.2 * 0.2 * 0.2 = 0.008$
-1	$1 - 0.001 - 0.006 - 0.008 - 0.008 = 0.977$

$$\text{Expectation } [x] = \sum_n x P(X=n)$$

$$= 19 * 0.001 + 14 * 0.006 + 9 * 0.008 + 4 * 0.008$$

$$= 1 * 0.977$$

$$E[x] = -0.77$$

$$\text{Variance } [x] = E[(x-\mu)^2]$$

$$\therefore E[f(n)] = \sum_n f(n) \cdot P(X=n)$$

$$= \sum_n (n-\mu)^2 P(X=n)$$

$$= \sum_n (n^2 + \mu^2 - 2\mu n) P(X=n)$$

$$= \sum_n n^2 P(X=n) + \mu^2 \sum_n P(X=n) - 2\mu \sum_n n P(X=n)$$

$$= E(x^2) - E^2(x)$$

X	Prob. ( $X=n$ )	$X-\mu$	$(X-\mu)^2$
19	$0.1 * 0.1 * 0.1 = 0.001$	$19 - 0.77$	$390.85$
14	$0.1 * 0.1 * 0.2 * 3 = 0.006$	$14 - 0.77$	$218.15$

9	$0.2 * 0.2 * 0.2 = 0.008$	9.77	95.45
4	$0.2 * 0.2 * 0.2 = 0.008$	4.77	22.75
-1	$1 - \text{AU} = 0.977$	-0.23	0.05

$$\begin{aligned}\text{var} &= 0.001 * 390.85 + 0.006 * 218.15 + 0.008 * 95.45 \\ &\quad + 0.008 * 22.75 + 0.977 * 0.05 \\ &= 2.69\end{aligned}$$

$$\sigma = \sqrt{2.69} = 1.64$$

- In general we lose 0.77 vs by gaussian dist.

$\therefore$  Our win range lie in  $[-0.77 - 3*1.64, -0.77 + 3*1.64]$

- Price money for the lottery is made 5 times & fee is doubled.

$Y \quad P(X)$

98	0.001	Original win - fee = $X$
73	0.006	Original win = $X + \text{fee} = (X+1)$
48	0.008	
23	0.008	$Y = 5(X+1) - 2$
-2	0.977	$Y = 5X + 3$



Linear dependency

- If  $Y = ax + b$ ;  $E(x)$  is known

$$\text{then } E[Y] = aE[X] + b$$

$$\text{var}[Y] = a^2 \text{var}[x]$$

$$\text{Derive} \Rightarrow E[X] = \sum_n f(n) \cdot P(X=n)$$

$$E[ax+b] = \sum_n (an+b) \cdot P(X=n)$$

$$\begin{aligned}
 &= \sum_n a \cdot n \cdot P(X=n) + \sum_n b \cdot P(X=n) \\
 &= a \sum_n n P(X=n) + b \\
 &= a E[n] + b
 \end{aligned}$$

$$\begin{aligned}
 \text{Var}(ax+b) &= E[(an+b)^2] - (E[an+b])^2 \\
 &= E[a^2n^2 + b^2 + 2abn] - (aE[n]+b)^2 \\
 &= a^2 E[x^2] + 2ab E[x] + b^2 - a^2 E^2[x] - 2ab E[x] - b^2 \\
 &= a^2 (E[x^2] - E^2[x]) \\
 &= a^2 \text{Var}[x]
 \end{aligned}$$

Note:  $2X$  is diff from  $X+X \Rightarrow 2X \Rightarrow$  price money doubled  
 $X+X$  is 2 lottery ticket purchase

$2X \rightarrow$ linear dependency			$X+X \rightarrow$ independent observations		
$X$	$15$	$-2$	$X$	$P(n)$	
$P(X=n)$	$0.4$	$0.6$	$30$	$0.16$	
$Y = 2X$	$30$	$-4$	$13$	$0.48$	
$P(Y=y)$	$0.4$	$0.6$	$-4$	$0.36$	

both differ

$$- E(x+x) = 2E(x)$$

$$\text{Var}[x+x] = 2\text{Var}[y]$$

$$- E(2x) = 2E(x)$$

$$\text{Var}[2x] = 4\text{Var}[y]$$

Derive:  $E(ax+by) = aE[x]+bE[y]$

$$\text{Var}[ax+by] = a^2 \text{Var}[x]+b^2 \text{Var}[y]$$

$$E[ax-by] = aE[x]-bE[y]$$

$$\text{Var}[ax-by] = a^2 \text{Var}[x]+b^2 \text{Var}[y]$$

For independent var

- Chap. 5 Head first statistics

Expectation & Variance of geometric distribution

P = Prob. win      r = trials

r = n to success, q = lose prob

$$E(x) = \sum_{r=1}^{\infty} r P(x=r)$$

$$= p \sum_{k=1}^{\infty} (k-1) k q^{k-1} = p \cdot \frac{d}{dq} \left( \sum_{k=1}^{\infty} (k-1) q^k \right)$$

$$E(x) = 1 \cdot p + 2 \cdot q p + 3 q^2 p + 4 q^3 p + \dots$$

$$= p \frac{d}{dq} \left( q^2 \sum_{k=1}^{\infty} (k-1) q^{k-1} \right) = p \frac{d}{dq} \left( q^2 \sum_{k=1}^{\infty} (k-1) q^{k-2} \right)$$

$$r E(x) = 1 \cdot q p + 2 q^2 p + 3 q^3 p + \dots$$

$$= p \frac{d}{dq} \left( q^2 \frac{d}{dq} \left( \sum_{k=1}^{\infty} q^{k-1} \right) \right) = p \frac{d}{dq} \left( q^2 \frac{d}{dq} \left( \sum_{k=1}^{\infty} q^k \right) \right)$$

$$(1-q) E(x) = p + p q + p q^2 + \dots + p q^k + \dots$$

$$= p \frac{d}{dq} \left( q^2 \frac{d}{dq} \left( \frac{1}{1-q} - 1 \right) \right) = p \frac{d}{dq} \left( q^2 \left( \frac{-1}{(1-q)^2} \right) \right)$$

$$P(T=r) = p [1 + q + q^2 + \dots + q^{r-1}]$$

$$= p \frac{d}{dq} \left( \frac{-q^r}{(1-q)^2} \right) = p \cdot \frac{2q}{(q-1)^3} = p \cdot \frac{2q}{p^3} = \frac{2q}{p^2}$$

$$E(r) = \frac{1}{1-q} = \frac{1}{p}$$

$$\therefore \text{Var}[r] = \frac{2(1-p)}{p^2} + \frac{1}{p} - \frac{1}{p^2}$$

$$\text{Var}[n] = E(x^2) - (E[x])^2$$

$$= 1-p/p^2$$

$$= E[x^2] - \frac{1}{p^2}$$

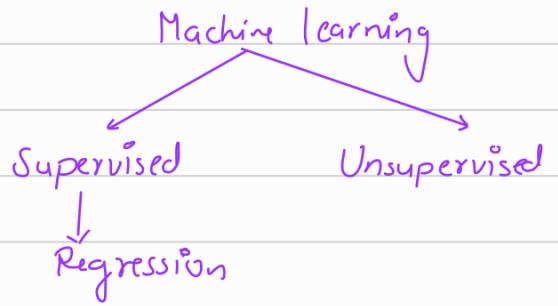
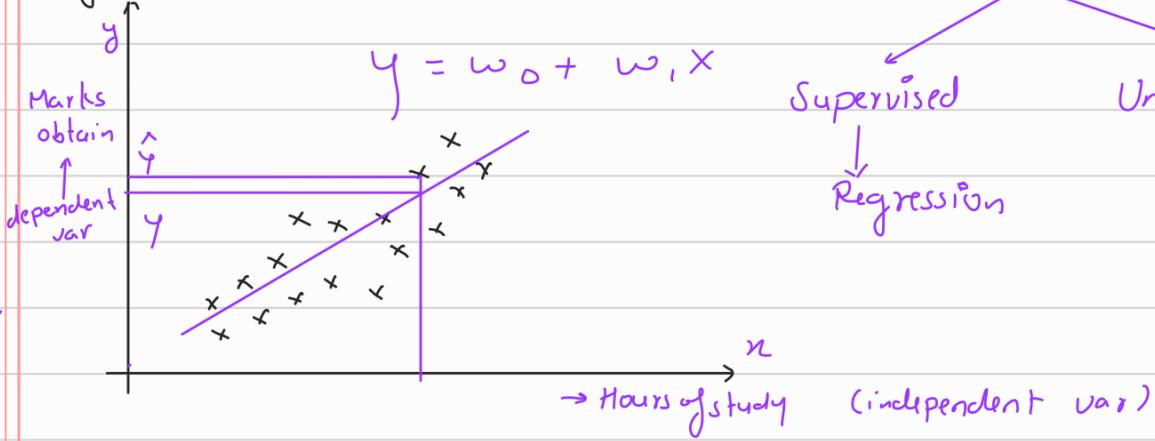
$$= q/p^2$$

$$= E[x(x-1)] + E[x] - (E[n])^2$$

$$E[x(x-1)] = \sum_{k=1}^{\infty} k(k-1) P(x=k)$$

$$= \sum_{k=1}^{\infty} k(k-1) p \cdot q^{k-1}$$

## Regression →



let Assume sleep also affect then

$$Y = w_1 \cdot X_1 + w_2 \cdot X_2 + w_0 \quad (\text{Multiple linear relation})$$

polynomial reg → let 1 ind. var but 2<sup>nd</sup> order relationships

$$Y = w_2 \cdot X^2 + w_1 \cdot X + w_0 \quad (\text{reg with 2<sup>nd</sup> order poly})$$

let 2 ind. var →

$$Y = w_0 + w_1 \cdot X_1 + w_2 \cdot X_2 + w_3 \cdot X_1^2 + w_4 \cdot X_2^2 + w_5 \cdot X_1 \cdot X_2$$

$\hat{Y}$	$X_1$	$X_2$	$\hat{Y} - Y$	$\frac{\text{sqv}}{n}$	$E = \frac{1}{2} \sum \text{sqv}$	$\text{let } Y = 2n$	$\gamma E$	$\propto w$
2	1	1	1	$\frac{\text{sqv}}{n}$	1	2	0	0
4	2	2	2	$\frac{\text{sqv}}{n}$	4	4	0	0
6	3	3	3	$\frac{\text{sqv}}{n}$	9	6	0	0
8	4	4	4	$\frac{\text{sqv}}{n}$	16	8	0	0

$$E = 15$$

this need to  
min.

$\therefore E=0$   
 $\downarrow \min$

- In regression we minimize least squares errors

$$E = \frac{1}{2} \sum_{i=1}^N y_i^2 = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

$$= \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - w_0 - w_1 n_i)^2$$

↓  
minimize

$$\frac{\partial E}{\partial w_0} = \frac{1}{2} \times 2 \sum_{i=1}^N (\hat{y}_i - w_0 - w_1 n_i) (-1) \Rightarrow 0$$

$$= \sum_{i=1}^N \hat{y}_i - w_0 N - w_1 \sum_{i=1}^N n_i = 0$$

$$= w_0 N + w_1 \sum_{i=1}^N n_i = \sum_{i=1}^N \hat{y}_i - \textcircled{1}$$

$$\frac{\partial E}{\partial w_1} = \frac{1}{2} \times 2 \sum_{i=1}^N (\hat{y}_i - w_0 - w_1 n_i) (-n_i) \Rightarrow 0$$

$$= \sum_{i=1}^N \hat{y}_i n_i - w_0 \sum n_i - w_1 \sum_{i=1}^N n_i^2 = 0 - \textcircled{2}$$

Use L, 2 & data, substitute it  $\Rightarrow$

$y$	$x$	$n^2$
2	1	1
4	2	4
6	3	9
8	4	16
20	10	30

$$\textcircled{1} \Rightarrow w_0(4) + w_1(10) = 20$$

$$\textcircled{2} \Rightarrow w_0(10) + w_1(30) = 60$$

$w_0 = 0 \quad w_1 = 2 \quad \Rightarrow \text{solution optimal}$

<u>Ex-</u>	$\hat{y}$	$n$	Fit in the following poly.
	4	1	(i) $\hat{y} = w_0 + w_1 n$
	8	2	(ii) $\hat{y} = w_0 + w_1 n + w_2 n^2$
	13	3	
	32	4	

(i) ①  $\Rightarrow w_0 \cdot 4 + w_1 \cdot 10 = 57$   
 ②  $\Rightarrow 187 = w_0 \cdot 10 + w_1 \cdot 30$

$$19w_0 + 30w_1 = 171$$

$$10w_0 + 30w_1 = 187$$

$$2w_0 = -16$$

$$w_0 = -8 \Rightarrow w_1 = 57 - 4(-8)$$

$$= 57 + 32 = 89$$

↗ n v  
↑ change

(ii) Solve diff.

Order = curve + 1

$\hat{y}$	$x$	$\hat{y} - y$
2	1	$2 - w$
5	2	$5 - 2w$
5	3	$5 - 3w$
7	4	$7 - 4w$
11	5	$11 - 5w$

$y = wn \rightarrow \text{Given}$

$$\frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2 = \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - wn)^2$$

$$\frac{dE}{dw} = \frac{1}{2} \times \sum_{i=1}^N -2(\hat{y}_i - wn_i)(-n_i) = 0$$

$$= \sum_{i=1}^N (\hat{y}_i - wn_i)n_i = 0$$

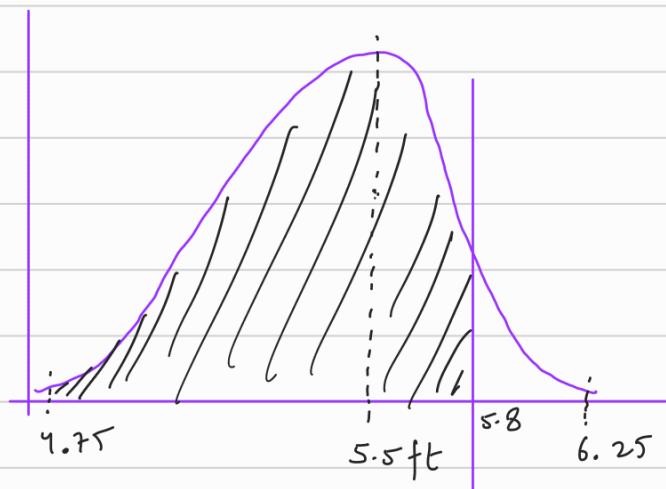
$$= \sum_{i=1}^N y_i n_i - w \sum_{i=1}^N n_i^2 = 0$$

$$= 110 - \omega 55 = 0$$

$$\Rightarrow \omega = 2$$

$$\therefore y = 2n$$

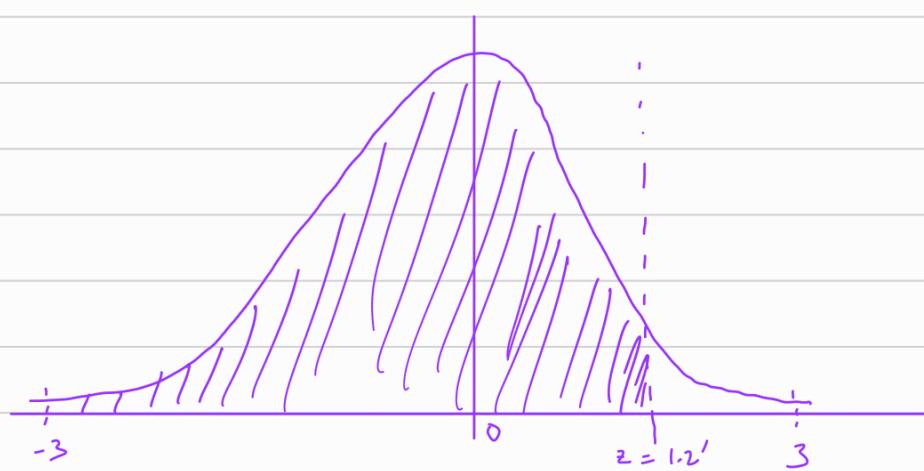
- Q. The height of students of class follow normal distribution with avg height of 5.5 ft and variance of  $1/16$  ft find out prob that student will have height less than 5.8 ft?



$$\sigma^2 = 1/16$$

$$\sigma = 1/4$$

Standard,  $z = \frac{x-\mu}{\sigma}$

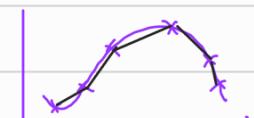


$$z = \frac{5.8' - 5.5'}{0.25} = 1.2'$$

$$\text{Ans} = z\text{-table}(1.2') =$$

SEE IN PYTHON

Curve fitting → Overfitting



taken deg. > req. deg.

Underfitting.



req. poly degree > taken poly

- To resolve both we use regularization

Regression with regularization →

- Here we modify error function to incorporate regularization term :-

$$\text{We need to minimise} \Rightarrow \frac{1}{2} \sum_{i=1}^N (\hat{y}_i - y_i)^2 + C \cdot \frac{1}{2} w^T w$$

↓

Regularization parameter

$C \in [2^{-20}, 2^{-18}, \dots, 2^{50}]$

$\begin{bmatrix} w_1, w_2 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \end{bmatrix}$  regularization

$$= \frac{1}{2} (w_1^2 + w_2^2)$$

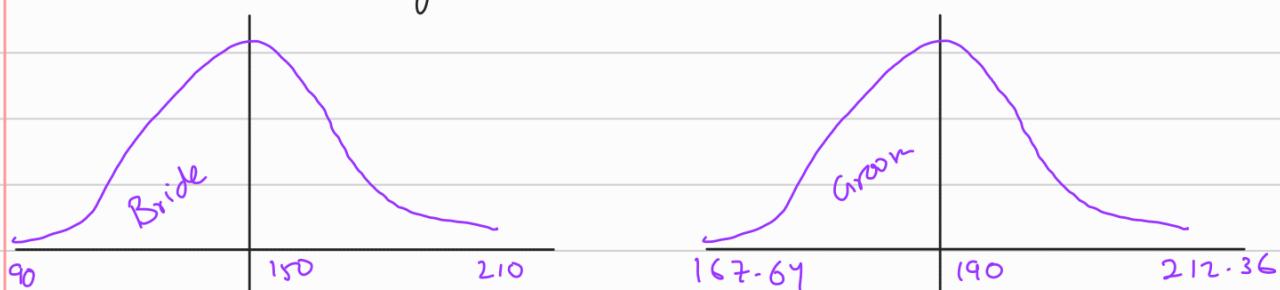
- In regularization we need to tune C & d (check all poss. for best result)
- KEEL DATASET Repo. → regression prob.
- Find 3-4 reg. prob. & use suitable tool box & solve these problems  
find optimal value of degree of polynomial & regularization polynomial

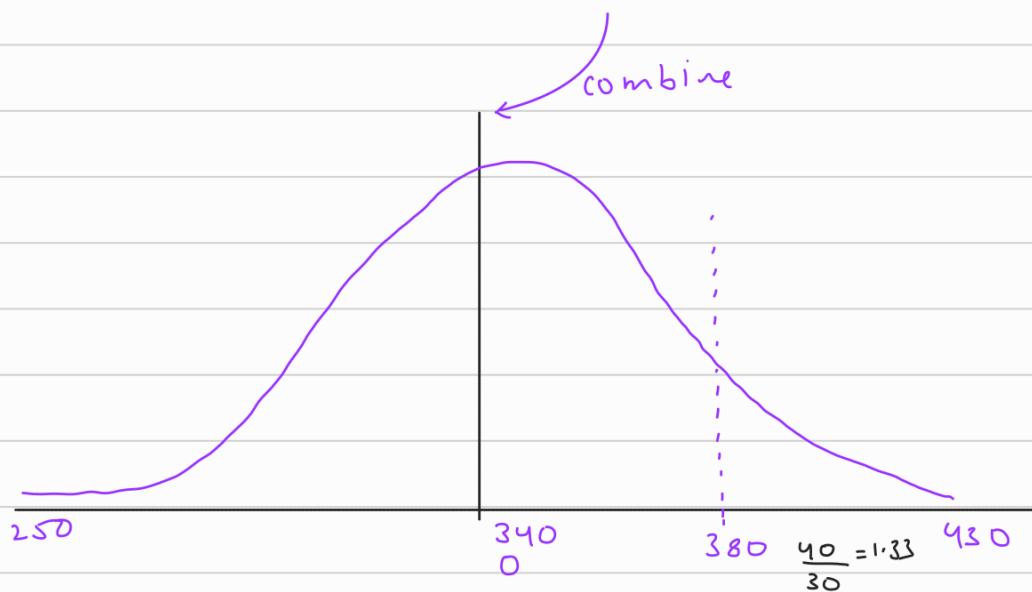
Chap-9 Beyond normal distr. →

Q.



- Q. Roller Coaster ride where  $X$  rep. weight of bride  $X \sim N(150, 400)$   
& weight of groom is given by  $Y \sim N(190, 500)$ , they can ride if their combined weight is  $\leq 380$





- When we sum up 2 distribution follow gaussian distribution it will also follow gaussian dist.

$$E(x+y) = E(x) + E(y)$$

$$\text{Var}[x+y] = \text{Var}[x] + \text{Var}[y]$$

$x, y$  - independent

$$\therefore (x+y) \sim N(340, 900)$$

$$Z = \frac{380 - 340}{30} = 1.33$$

$$\therefore \text{Area} = \text{Prob} = 0.908$$

- In town male height  $\sim N(71, 20.25)$ , & for women  $\sim N(64, 16)$  what is prob that man will be at least 5 inches taller than women

$$X : N(71, 20.25)$$

$$X - Y \geq 5$$

$$Y : N(64, 16)$$

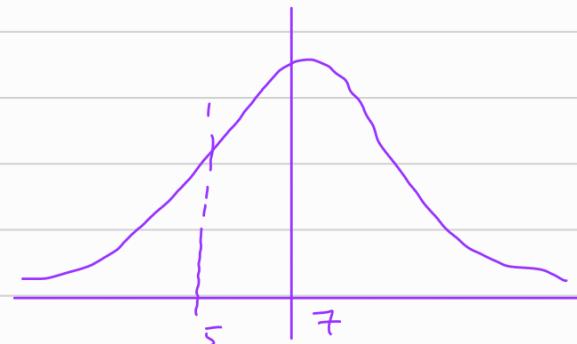
$x, y \in$  independent

$$E(X-Y) = E(X) - E(Y) = E_N$$

$$\text{Var}[X-Y] = \text{Var}[X] + \text{Var}[Y] = V_N$$

$$E_N = 7, V_N = 36.25$$

$$\sigma_N = 6.02$$



$$P(<5) = Z \left( \frac{5-7}{6.02} \right) = Z \left( \frac{-2}{6.02} \right) = Z(-0.33) = 0.37070$$

$$\therefore P(>=5) = 1 - 0.37070 \\ \approx 0.63$$

Till chapter 9  
Read first stat.

### Binomial distribution →

- 'n' - Toss head success
- $P(X=r) = {}^n C_r p^r q^{n-r}$

Ques. paper 50 Q → 4 possible ans (A, B, C, D) one of them is correct prob  $\leq 30$  correct?

$$p = \frac{1}{4}, q = \frac{3}{4}, n = 50$$

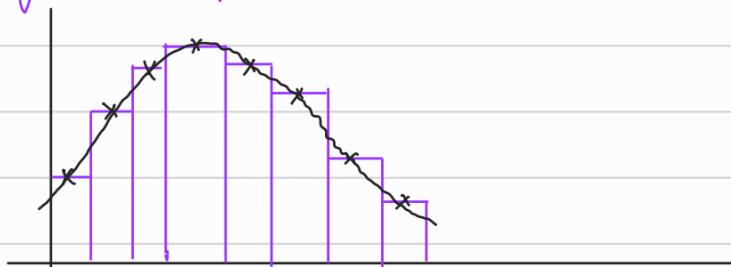
$$P(X \leq 30) = \sum_{r=0}^{30} {}^n C_r p^r q^{n-r}$$

( )  
→ Calculating

∴ If  $np \leq 10$  we plot a histogram for 'i' we plot bar from  $i-0.5$  to  $i+0.5$

→ but if  $np > 10$  then this histogram tends to follow gaussian distribution.

- Then we normalize freq. of histogram by dividing freq./Total freq.
- Then we apply regression to find the curve



- let by reg we get pdf if  $\int_{-\infty}^{\infty} \text{pdf} \cdot dn = P$  divide it by

$$P \therefore \text{pdf} = \frac{1}{P} (\text{pdf})$$

Else we can use progen window estimator -

- let 3 student height 5', 5.5', 5.6' & i want to find overall pdf
- we will first find standard deviation ' $\sigma$ '
- we will take diff function

$$\text{pdf}_{\text{Progen}} = \frac{1}{3} (GF(5', \sigma) + GF(5.5', \sigma) + GF(5.6', \sigma))$$

- The shape of window can be anything

$$E[X] = np \\ \text{Var}[X] = npq$$

$X$	0	1
$P(X)$	$q$	$p$

$$E[X]_{\text{single trial}} = \sum_{x=0} x P(x=n) = 0 \cdot q + 1 \cdot p = p$$

$$E[X_1 + X_2 + \dots + X_n] = n p \quad \because \text{trials are independent}$$

$$E[X] = np$$

$$\text{Var}[X_1 + \dots + X_n] = n \text{Var}[X]$$

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

single trial

$$\begin{aligned}
 &= \sum x^2 P(x=n) - p^2 \\
 &= 0 \cdot q + 1 \cdot p - p^2 \\
 &= p - p^2 = p(1-p) = pq
 \end{aligned}$$

$$\therefore \text{Var}[x] = Pq$$

for 1 trial

$$\text{So. } \text{Var}_n[x] = nPq$$

Q. There are 20 Question the prob of getting 4 or less question correct (single MCQ ans 2 option)

$$P(X=4) = {}^{20}C_4 \left(\frac{1}{2}\right)^4 \left(\frac{1}{2}\right)^{16}$$

$$P(X=3) = {}^{20}C_3 \left(\frac{1}{2}\right)^{20}$$

$$P(X=2) = {}^{20}C_2 \left(\frac{1}{2}\right)^{20}$$

$$P(X=1) = {}^{20}C_1 \left(\frac{1}{2}\right)^{20}$$

$$P(X=0) = {}^{20}C_0 \left(\frac{1}{2}\right)^{20}$$

$$\Sigma = \frac{1}{2^{20}} \left[ {}^{20}C_0 + {}^{20}C_1 + {}^{20}C_2 + {}^{20}C_3 + {}^{20}C_4 \right]$$

$$= \frac{1}{2^{20}} \left[ 1 + 20 + 190 + 1140 + 4845 \right]$$

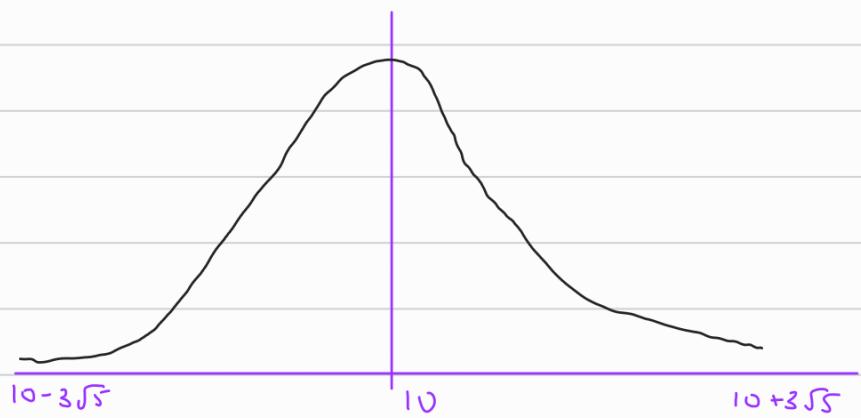
$$= \frac{1}{2^{20}} \cdot (6196) = 0.0059$$

- Expectation is mean.

Transform

$$\text{Exp}[X] = np = 20 \cdot \frac{1}{2} = 10 \geq 10 \therefore \text{Gaussian distribution}$$

$$\mu = 10, \quad \text{Var} = npq = 5 \Rightarrow \text{SD} = \sqrt{5}$$



$X = 4$  transform

$$Z = \frac{4 - 10}{\sqrt{5}} = -2.68$$

$$\text{Area } P(X \leq 4) = 0.0037$$

Wrong

we need to apply distribution correction

for histogram goes till 10.5

$\therefore X = 4.5$  transform

$$Z = \frac{4.5 - 10}{\sqrt{5}} = \frac{-5.5}{\sqrt{5}} = -1.1 * \sqrt{5} = -2.46$$

$$\text{Area} = \text{Area}(Z) = 0.0069 \approx 0.0059 \therefore \text{Right}$$

## Questions

**1.98 Lifetimes of Flashlight Batteries.** Two different options are under consideration for comparing the lifetimes of four brands of flashlight battery, using 20 flashlights.

- One option is to randomly divide 20 flashlights into four groups of 5 flashlights each and then randomly assign each group to use a different brand of battery. Would this statistical design be a completely randomized design or a randomized block design? Explain your answer.
- Another option is to use 20 flashlights—five different brands of 4 flashlights each—and randomly assign the 4 flashlights of each brand to use a different brand of battery. Would this statistical design be a completely randomized design or a randomized block design? Explain your answer.

- Complete random no generator when 3 same unit (actus) we will divide them by random no without replacement
- randomize block design when 3 more than 1 species say out of 16 ] (8, 8), 2 species then we divide them equally both species

	T1	T2	T3	T4
2 SPL	2	2	2	2
2 SP.2	2	2	2	2

**1.91 Treating Heart Failure.** In the journal article "Cardiac Resynchronization Therapy with or without an Implantable Defibrillator in Advanced Chronic Heart Failure" (*New England Journal of Medicine*, Vol. 350, pp. 2140–2150), M. Bristow et al. reported the results of a study of methods for treating patients who had advanced heart failure due to ischemic or nonischemic cardiomyopathies. A total of 1520 patients were randomly assigned in a 1:2:2 ratio to receive optimal pharmacologic therapy alone or in combination with either a pacemaker or a pacemaker-defibrillator combination. The patients were then observed until they died or were hospitalized for any cause.

- How many treatments were there?
  - Which group would be considered the control group?
  - How many treatment groups were there? Which treatments did they receive?
  - How many patients were in each of the three groups studied?
  - Explain how a table of random numbers or a random-number generator could be used to divide the patients into the three groups.
- In Exercises 1.92–1.97 ...

- optimal therapy
  - Pacemaker optimal therapy
  - Pacemaker-defibrillator optimal therapy
- Optimal therapy
- $\frac{1520}{5}$ ,  $\frac{1520 \times 2}{5}$ ,  $\frac{1520 \times 2}{5}$

3 treatment

Jisko  
desired  
treatment  
denge

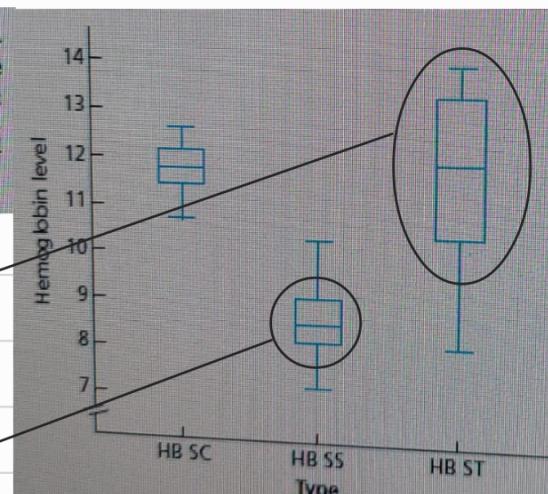
d.

e.

**3.181 Sickle Cell Disease.** A study published by E. Anionwu et al. in the *British Medical Journal* (Vol. 282, pp. 283–286) examined the steady-state hemoglobin levels of patients with three different types of sickle cell disease: HB SC, HB SS, and HB ST. Use the following boxplots to compare the hemoglobin levels for the three groups of patients, paying special attention to center and variation.

High variation

low sugar level



**23. Millionaires.** Refer to Problem 20. The ages of the 36 millionaires sampled are arranged in increasing order in the following table.

31	38	39	39	42	42	45	47	48
48	48	52	52	53	54	55	57	59
60	61	64	64	66	66	67	68	68
69	71	71	74	75	77	79	79	79

- a. Determine the quartiles for the data.
- b. Obtain and interpret the interquartile range.
- c. Find and interpret the five-number summary.
- d. Calculate the lower and upper limits.
- e. Identify potential outliers, if any.
- f. Construct and interpret a boxplot.

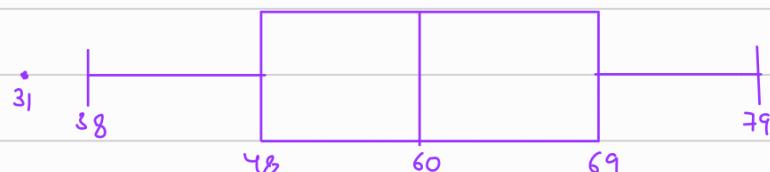
Outlier represented by dot in boxplot isko ignore kro aur median nikalo

$$(M) \text{ median} = 60 \quad (\text{35 element})$$

$$\text{median } (1 \rightarrow M) = 48 \quad (\text{17 element})$$

$$\text{median } (M+1 \rightarrow \text{last}) = 69 \quad (\text{17 element})$$

} Quartile



Outlier : 31.

$$\text{Interquartile range: } 69 - 48 = 21$$

$$L = 38, U = 79.$$

3.195 Fill in the following blanks.

- A standardized variable always has mean \_\_\_\_\_ and standard deviation \_\_\_\_\_.
  - The z-score corresponding to an observed value of a variable tells you \_\_\_\_\_.
  - A positive z-score indicates that the observation is \_\_\_\_\_ the mean, whereas a negative z-score indicates that the observation is \_\_\_\_\_ the mean.
- 3.196 Identify the statistic that is used to estimate
- a population mean.
  - a population standard deviation.

3.216 Low-Birth-Weight Hospital Stays. Data on low-birth-weight babies were collected over a 2-year period by 14 participating centers of the National Institute of Child Health and Human Development Neonatal Research Network. Results were reported by J. Lemons et al. in the on-line paper "Very Low Birth Weight Outcomes of the National Institute of Child Health and Human Development Neonatal Research Network" (*Pediatrics*, Vol. 107, No. 1, p. e1). For the 1084 surviving babies whose birth weights were 751–1000 grams, the average length of stay in the hospital was 86 days, although one center had an average of 66 days and another had an average of 108 days.

- Can the mean lengths of stay be considered population means? Explain your answer.
- Assuming that the population standard deviation is 12 days, determine the z-score for a baby's length of stay of 86 days at the center where the mean was 66 days.
- Assuming that the population standard deviation is 12 days, determine the z-score for a baby's length of stay of 86 days at the center where the mean was 108 days.
- What can you conclude from parts (b) and (c) about an infant with a length of stay equal to the mean at all centers if that infant was born at a center with a mean of 66 days? mean of 108 days?

Q.  $B \sim N(65, 16)$ ,  $G \sim N(60, 9)$

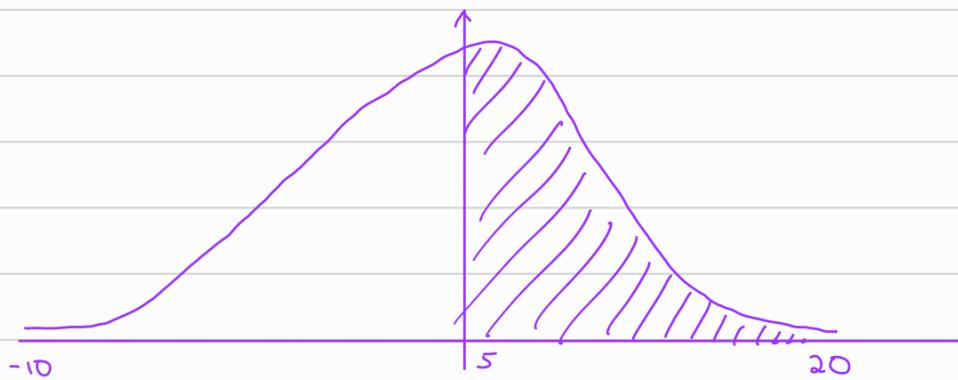
Prob. that a random boy chosen is 5 or more inch taller than a randomly chosen girl

B:  $\mu = 65$ ,  $\sigma = 4$

G:  $\mu = 60$ ,  $\sigma = 3$

$$E(N) = E_B - E_G = 65 - 60 = 5$$

$$\text{Var}(N) = \text{Var}_B + \text{Var}_G = 16 + 9 = 25 \Rightarrow \sigma_N = 5$$



$$Z = \frac{5-5}{5} = 0$$

$$\therefore \text{Area} = \text{Prob} = \text{z-score}(0) = 0.5$$

Q. Fit the poly.  $y = w_0 + w_1 x_1$  for the following data. Use ridge regularization take  $C=10$ . Also do it without regularization

$$(y, x) = (3, 1), (4, 2), (5, 3), (4, 4), (5, 5), (7, 6)$$

- Without regularization -

$$\Rightarrow \sum_{i=1}^N \hat{y}_i n_i - w_0 \sum n_i - w_1 \sum_{i=1}^N n_i^2 = 0 \quad \textcircled{2}$$

$$\Rightarrow w_0 N + w_1 \sum_{i=1}^N n_i = \sum_{i=1}^N \hat{y}_i \quad \textcircled{1}$$

$$\textcircled{2} \Rightarrow (3+8+15+16+25+42) - w_0(21) - w_1(91) = 0$$

$$\begin{array}{c} \hat{y} \\ 3 \\ 4 \\ 5 \\ 4 \\ 5 \\ 7 \end{array} \quad \begin{array}{c} x \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \quad 109 = 21w_0 + 91w_1$$

$$\textcircled{1} \Rightarrow w_0(6) + w_1(21) = 28$$

$$\begin{array}{c} \hat{y} \\ 3 \\ 4 \\ 5 \\ 4 \\ 5 \\ 7 \end{array} \quad \begin{array}{c} x \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \quad w_0 = \frac{28 - 21w_1}{6}$$

$$109 = 21\left(\frac{28 - 21w_1}{6}\right) + 91w_1$$

$$218 = 196 - 147w_1 + 182w_1$$

$$w_1 = 22/35$$

- With regularization -

$$y = w_0 + w_1 n \quad , \quad C = 10$$

$$\begin{array}{c} \hat{y} \\ 3 \\ 4 \\ 5 \\ 4 \\ 5 \\ 7 \end{array} \quad \begin{array}{c} x \\ 1 \\ 2 \\ 3 \\ 4 \\ 5 \\ 6 \end{array} \quad \begin{array}{c} y \\ w_0 + w_1 \\ w_0 + 2w_1 \\ w_0 + 3w_1 \\ w_0 + 4w_1 \\ w_0 + 5w_1 \\ w_0 + 6w_1 \end{array} \quad \begin{array}{c} y - \hat{y} \\ | \\ | \\ | \\ | \\ | \\ | \end{array}$$

$$3 \quad 1 \quad w_0 + w_1$$

$$4 \quad 2 \quad w_0 + 2w_1$$

$$5 \quad 3 \quad w_0 + 3w_1$$

$$4 \quad 4 \quad w_0 + 4w_1$$

$$5 \quad 5 \quad w_0 + 5w_1$$

$$7 \quad 6 \quad w_0 + 6w_1$$

$$C. \frac{1}{2} \sum (y - q)^2 + \frac{1}{2} \cdot \left( \begin{bmatrix} w_0 \\ w_1 \end{bmatrix} \begin{bmatrix} w_0 & w_1 \end{bmatrix} \right) \text{ to be minim}$$

$$(w_0 + w_1 - 3)^2 + (w_0 + 2w_1 - 4)^2 + (w_0 + 3w_1 - 5)^2 + (w_0 + 4w_1 - 4)^2 + (w_0 + 5w_1 - 5)^2 + (w_0 + 6w_1 - 7)^2 + \frac{1}{10} (w_0^2 + w_1^2) \text{ to be minim}$$

$$\begin{aligned} & (w_0 + w_1)^2 + 9 - 6w_0 - 6w_1, \\ & + (w_0 + 2w_1)^2 + 16 - 8w_0 - 16w_1, \\ & + (w_0 + 3w_1)^2 + 25 - 10w_0 - 30w_1, \\ & + (w_0 + 4w_1)^2 + 16 - 8w_0 - 16w_1, \\ & + (w_0 + 5w_1)^2 + 25 - 10w_0 - 50w_1, \\ & + (w_0 + 6w_1)^2 + 49 - 14w_0 - 84w_1, \\ & + w_0^2 + w_1^2 \quad \text{to min} \end{aligned}$$

$$\Rightarrow (w_0 + w_1)^2 + (w_0 + 2w_1)^2 + (w_0 + 3w_1)^2 + (w_0 + 4w_1)^2 + (w_0 + 5w_1)^2 + (w_0 + 6w_1)^2 + \frac{1}{10} (w_0^2 + w_1^2) + 140 - 56w_0 - 202w_1 \text{ to be mi.}$$

$$\begin{aligned} \Rightarrow & w_0^2 + w_1^2 + 2w_0w_1 + w_0^2 + 4w_1^2 + 4w_0w_1 + w_0^2 + 9w_1^2 + 6w_0w_1, \\ & + w_0^2 + 16w_1^2 + 8w_0w_1 + w_0^2 + 25w_1^2 + 10w_0w_1 + w_0^2 + 36w_1^2 \\ & + 12w_0w_1 + \frac{1}{10}(w_0^2 + w_1^2) - 56w_0 - 202w_1, \end{aligned}$$

$$\frac{\delta E}{\delta w_0} = 0$$

$$\frac{\delta E}{\delta w_1} = 0$$

$$61w_0^2 + 91w_1^2 + 420w_0w_1 - 560w_0 - 2020w_1 \quad \text{min}$$

$$122w_0 + 420w_1 = 560 \Rightarrow 61w_0 + 210w_1 = 280$$

$$1822w_1 + 420w_0 = 2020 \Rightarrow 911w_1 + 210w_0 = 1010$$

$$w_0 = \frac{1010 - 911w_1}{210}$$

$$\frac{61(1010 - 911\omega_1)}{210} + 210\omega_1 = 280$$

$$61(1010 - 911\omega_1) + 44100\omega_1 = 280 \times 210$$

$$(44100 - 911 \times 61)\omega_1 = 280 \times 210 - 61 \times 1010$$

$$-11471\omega_1 = -2810$$

$$\omega_1 = 0.25$$

$$\omega_0 = 3.725$$

Q 4.72

$$\mu = 78, \sigma = 10$$

Q. 4.73

$$\mu = ?, \sigma = ?$$

$$z = \frac{70 - \mu}{\sigma} = -0.6$$

$$z = \frac{88 - \mu}{\sigma} = 1.4$$

$$70 - \mu = -0.6\sigma$$

$$88 - \mu = 1.4\sigma$$

$$18 = 2\sigma$$

$$\sigma = 9 \quad \therefore \mu = 70 + 0.6(9) \\ = 70 + 5.4 = 75.4$$

Q. 4.74

(a)  $z = -1.20$  to  $z = 2.40$

$$z(2.40) - z(-1.20) = 0.99180 - 0.11507 \\ = 0.87673$$

(b)  $z = 1.23$  to  $z = 1.87$

$$z(1.87) - z(1.23) = 0.96926 - 0.89065 \\ = 0.07861$$

(c)  $z = -2.35$  to  $z = -0.5$

$$z(-2.35) - z(-0.5) = 0.30854 - 0.00939 \\ = 0.29915$$

Q.4.75

- (a)  $z(-1.78) = 0.03754$
- (b)  $z(0.56) = 0.71226$
- (c)  $1 - z(-1.45) = 1 - 0.07353 = 0.92647$
- (d)  $1 - z(2.16) = 1 - 0.98461 = 0.01539$
- (e)  $z(1.53) - z(-0.80) = 0.93699 - 0.21186 = 0.72513$
- (f)  $z(-2.52) + 1 - z(1.83) = 0.00587 + 1 - 0.96638$   
 $= 0.03949$

Q.4.76

- (a)  $1 - z(-1.64) = 1 - 0.05050 = 0.94949$
- (b)  $z(1.96) - z(-1.96) =$
- (c)  $z(-1) + 1 - z(1) =$

Q.4.77

- (a)  $1 - 0.2266 = 0.7734 \Rightarrow z = 0.75$
- (b)  $0.0314 \Rightarrow z = -1.86$
- (c)  $z(z) - z(-0.25) = 0.5722$

$$z(z) - 0.40905 = 0.5722$$

$$\begin{aligned} z(z) &= 0.98125 \\ &= 2.08 \end{aligned}$$

Q.4.78

$$1 - z(z_1) = 0.84$$

$$z(z_1) = 0.16$$

$$z_1 = 0.99$$

4.79

$$\mu = 5, \sigma = 2$$

$$\begin{aligned} P(x > 8) &= 1 - P(x \leq 8) = 1 - z((8-5)/2) \\ &= 1 - z(1.5) \\ &= 1 - 0.93319 \\ &= 0.06681 \end{aligned}$$

Q.4.80

$$N = 300, \mu = 68, \sigma = 3$$

$$\begin{aligned} (i) (1 - P(x > 72)) \cdot N &= (1 - z((72 - 68)/3)) \cdot N \\ &= (1 - z(1.33)) \cdot 300 \\ &= (1 - 0.90824) \cdot 300 \\ &= [0.09176 \times 300] \\ &= [27.5] \\ &= 28 \end{aligned}$$

$$\begin{aligned}
 \text{(ii)} \quad P(X < 64) \cdot N &= [z((64-68)/3)] \cdot N^7 \\
 &= [0.09342 \times 300] \\
 &= [28.0267] = 29
 \end{aligned}$$

$$\begin{aligned}
 \text{(iii)} \quad (P(X < +1) - P(X < -1))N &= [z\left(\frac{71-68}{3}\right) - z\left(\frac{65-68}{3}\right)] N \\
 &= [z(1) - z(-1)] N \\
 &= [0.84134 - 0.15866] 300 \\
 &= [204.8] = 205
 \end{aligned}$$

$$\text{(iv)} \quad P(X = 68) \cdot N = ?$$

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

$$P(68) = \frac{1}{\sqrt{2\pi} \times 3} e^{-\frac{(68-68)^2}{2 \times 9}} = \frac{1}{\sqrt{2\pi} \times 3} = 0.133$$

$$N(x) = P(68) \cdot N = 0.133 \times 300 = [39.9] = 40$$

Q. 4.81

$$\mu = 0.6140, \sigma = 0.0025$$

$$\begin{aligned}
 \text{(i)} \quad &[P(X < 0.618) - P(X < 0.616)] \cdot 100 \\
 \Rightarrow &[z(1.6) - z(-1.6)] \cdot 100 \\
 \Rightarrow &[0.94520 - 0.05480] \cdot 100 \\
 \Rightarrow &89.04\%
 \end{aligned}$$

$$\begin{aligned}
 \text{(ii)} \quad [P(X > 0.617)] &= 1 - P(X < 0.617) \\
 &= 1 - P(1.2) \Rightarrow 11.507\%
 \end{aligned}$$

$$(iii) [P(X < 0.608)] \times 100 = Z(-2.4) \times 100 = 6.00820 \times 100 \\ = 0.820\%$$

$$(iv) P(X = 0.615) \times 100 = 100 \cdot \frac{1}{\sqrt{2\pi} \sigma} e^{-\frac{(X-\mu)^2}{2\sigma^2}} \\ = 100 \cdot \frac{1}{\sqrt{2\pi} \times 0.0025} \times e^{-\frac{(0.615-0)^2}{2(0.0025)^2}} \\ = \frac{100 \times 400}{\sqrt{2\pi} \times e^{0.08}} = 14710\%.$$

Q. 4.82

$$\mu = 72, \sigma = 9$$

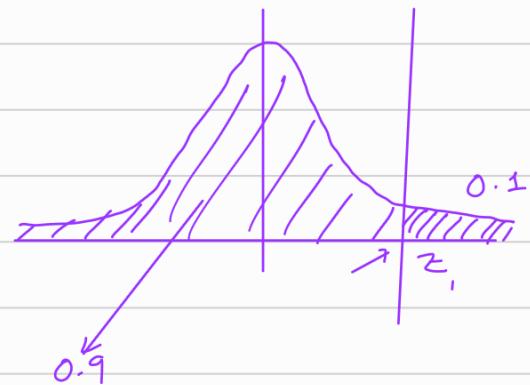
Top 10% = A

min score A

$$Z\left(\frac{z_1 - 72}{9}\right) = 0.9$$

$$z_1 - 72 = 1.29 \times 9$$

$$z_1 = 83.61$$



Q. 4.83

$$(a) P(X < -0.5) + 1 - P(X > 0.5)$$

$$0.30854 + 1 - 0.69146$$

$$0.61708$$

$$(b) P(X < 0.75)$$

$$0.77337$$

Poisson distribution  $\rightarrow$

$$P(X = r) = \frac{e^{-\lambda} \lambda^r}{r!}$$

Q. 4.84

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

(a)  $P(X < 2) - P(X < -2)$

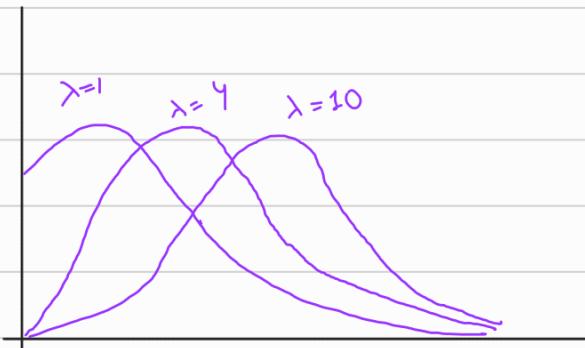
(b)  $P(X < -1.2) + 1 - P(X < 1.2)$

(c)  $1 - P(X < -1.5)$

Q. 4.85  $\omega P(a) - P(-a) = 0.75$

(b)  $P(-a) = 0.22$

Poisson dist. for diff  $\lambda$



$\lambda = 1$  Right skewed  
 $\lambda = 4$  less right skewed  
 $\lambda = 10$  Almost gaussian

In Binomial distribution  $np = \lambda$

( $n$  is very large  $p$  is very small)

$X \sim P(\lambda)$

$X$  follows poisson

$$E(X) = np$$

$$V[X] = npq \quad \because p \text{ is very small} \quad \therefore q \approx 1$$

$$\boxed{np \geq 10 \quad \& \quad npq \leq np}$$

We can transform  
poisson to  
normal distribu  
-ion.

$$\therefore E[X] = V[X] = \lambda$$

If  $\lambda \geq 10$  we can transform poisson distribution to normal distribution.

Q. A student need to take exam but has not done any rev. for it. he need to guess to ans to each ques. prob to guess right ans = 0.05  
30 Ques in exam what is prob he get 5 Q right Use poisson approximation to normal distribution?

$$n = 50, p = 0.05, X = 5$$

$$\lambda = n.p = 2.5$$

$$E(X = r) = \frac{e^{-\lambda} \lambda^r}{r!}$$

- Prove that

$$\lim_{n \rightarrow \infty} n c_r p^r (1-p)^{n-r}$$

$$\lim_{n \rightarrow \infty} \frac{n!}{r!(n-r)!} \left(\frac{\lambda}{n}\right)^r \left(1-\frac{\lambda}{n}\right)^{n-r}$$

$$\frac{\lambda^r}{r!} \lim_{n \rightarrow \infty} \frac{n!}{(n-r)!} \left(\frac{1}{n}\right)^r \left(1-\frac{\lambda}{n}\right)^n$$

$$n(n-1)(n-2) \dots (n-r+1)$$

$n \cdot n \cdot n \dots r \text{ times}$

$$\Rightarrow 1$$

$$\frac{\lambda^r}{r!} \lim_{n \rightarrow \infty} \left(1 + \frac{1}{\frac{-\lambda}{\lambda}}\right)^{\left(\frac{-\lambda}{\lambda}\right)(-r)} \Rightarrow \frac{e^{-\lambda} \lambda^r}{r!} \quad \because a^{mn} = (a^m)^n$$

$$\therefore \lim_{n \rightarrow \infty} \left(1 + \frac{1}{n}\right)^n = e$$

- Expected value of binomial distribution (long way)

$$E(K) = \sum_{k=0}^n k \cdot n c_k p^k q^{n-k} = np \quad (\text{Proof})$$

$$= \sum_{k=1}^n k \cdot n c_k p^k q^{n-k}$$

$$= \sum_{k=1}^n k \cdot \frac{n}{k} \cdot n-1 c_{k-1} p^k q^{n-k}$$

$$= \sum_{k=1}^n n \cdot n-1 c_{k-1} p^{k-1} q^{n-k} \cdot p$$

$$= np \sum_{k=1}^n n-1 c_{k-1} p^{k-1} q^{n-k}$$

$$= np (p+q)^{n-1}$$

$$= np$$

Proof :-

- Exp. & Var. of Geometric distribution
- Exp & var of binomial distribution
- Binomial to poisson distribution

mini

Imp topic :-

Boxplot

1 marks

Histogram

2 marks

Regression

2 marks

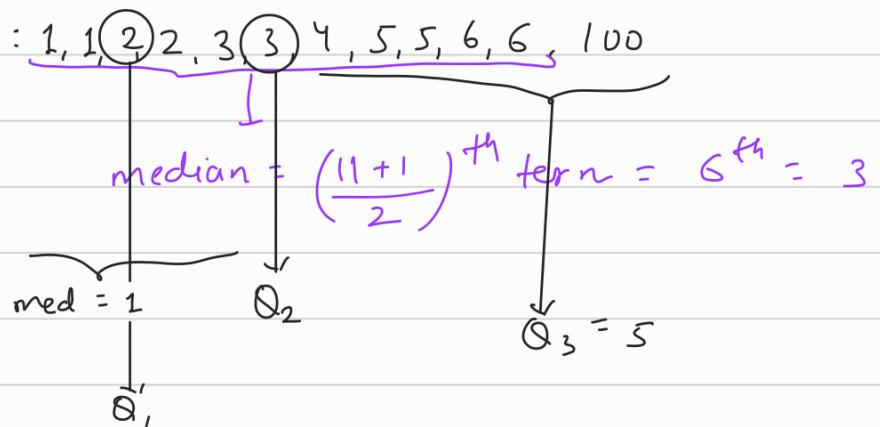
Normal distribution , binomial dist. poisson  
binomial to normal

3 marks

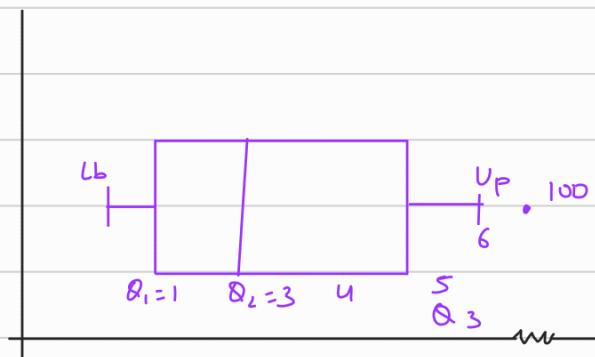
2 marks

Q. Draw Box plot : 1, 1, 2, 2, 3, 3, 4, 5, 5, 6, 6, 100

Outlier = 100



learn how to analyze  
boxplot



Q. Histogram

Age of mem of Clgss

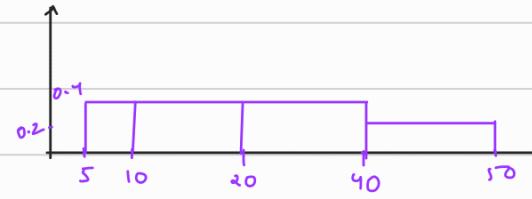
Int. width h:

5 - 10      2      5      0.4

10 - 20      4      10      0.4

20 - 40      8      20      0.4

40 - 50      2      10      0.2



Q. Regression →

Lasso regularization is also there but it is calculative (Used in ML)

Find out normal eqn

$$E = C \frac{1}{2} \sum_{i=1}^N (\hat{Y} - y_i)^2 + \frac{1}{2} (w_0^2 + w_1^2)$$

$$\text{where } \hat{y} = w_0 + w_1 n$$

$$E = \frac{1}{2} C \sum_{i=1}^N (\hat{y} - w_0 - w_1 n)^2 + \frac{1}{2} (w_0^2 + w_1^2)$$

$$\frac{\partial E}{\partial w_0} = 0, \quad \frac{\partial E}{\partial w_1} = 0$$

Regression (Matrix formulation) → (OLS using Wikipedia)  
using moore penrose

$$Y_{N \times 1} = X_{N \times m} \beta_{m \times 1} \quad \text{We need to find } \beta \quad \text{pseudo inverse}$$

$$\text{We need to minimize } \| Y - X\beta \|^2 = (X\beta - Y)_{1 \times N}^T (X\beta - Y)_{N \times 1}$$

$$L = (\beta^T X^T - Y^T)(X\beta - Y)$$

$$L = \beta^T X^T X \beta - \beta^T X^T Y - Y^T X \beta + Y^T Y$$

$$\frac{\partial L}{\partial \beta} = -2X^T Y + 2X^T X \beta = 0$$

$$(X^T X)^{-1} X^T = \text{Moore penrose pseudo inverse}$$

↳ it minimizes least square error.

With regularization

$$\beta = (X^T X + I/c)^{-1} X^T Y$$

$I$  = Identity of size  $m \times m$

$c$  = Reg. param

①  $\frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2 \Rightarrow (Y - X\beta)^T (X\beta - Y)$

②  $C \cdot \frac{1}{2} \sum_{i=1}^n (\hat{y}_i - y_i)^2 + \beta_0^2 + \beta_1^2 + \dots$

↓

$$\frac{C}{2} (X\beta - Y)^T (X\beta - Y) + \frac{1}{2} \beta \beta^T$$
$$\Rightarrow \beta = (X^T X + I/c)^{-1} X^T Y$$

A1 → Implement regularized regression without using any library function tune the value of regularization parameter  $C$  to obtain optimal result

A2 → submit the stepwise proof of regularization regression in matrix form

