

- 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

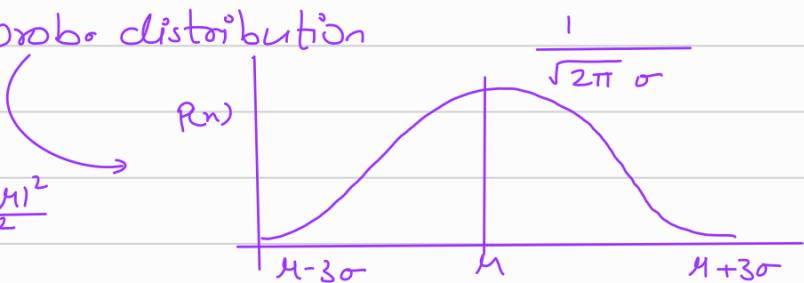
let $y = w_0 + w_1 n_1 + w_2 n_2 + w_3 n_3 + w_4 \cdot n_1^2 + w_5 n_2^2 + w_6 n_3^2$
 $+ w_7 n_1 n_2 + w_8 n_2 n_3 + w_9 n_1 n_3$

- 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 (σ) ?
- 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, σ . 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, σ

level of factor → C = 35, σ = 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

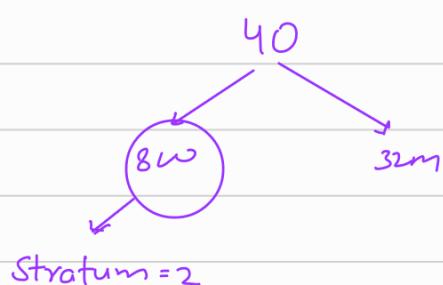
- Complete randomized design
- Randomized block design

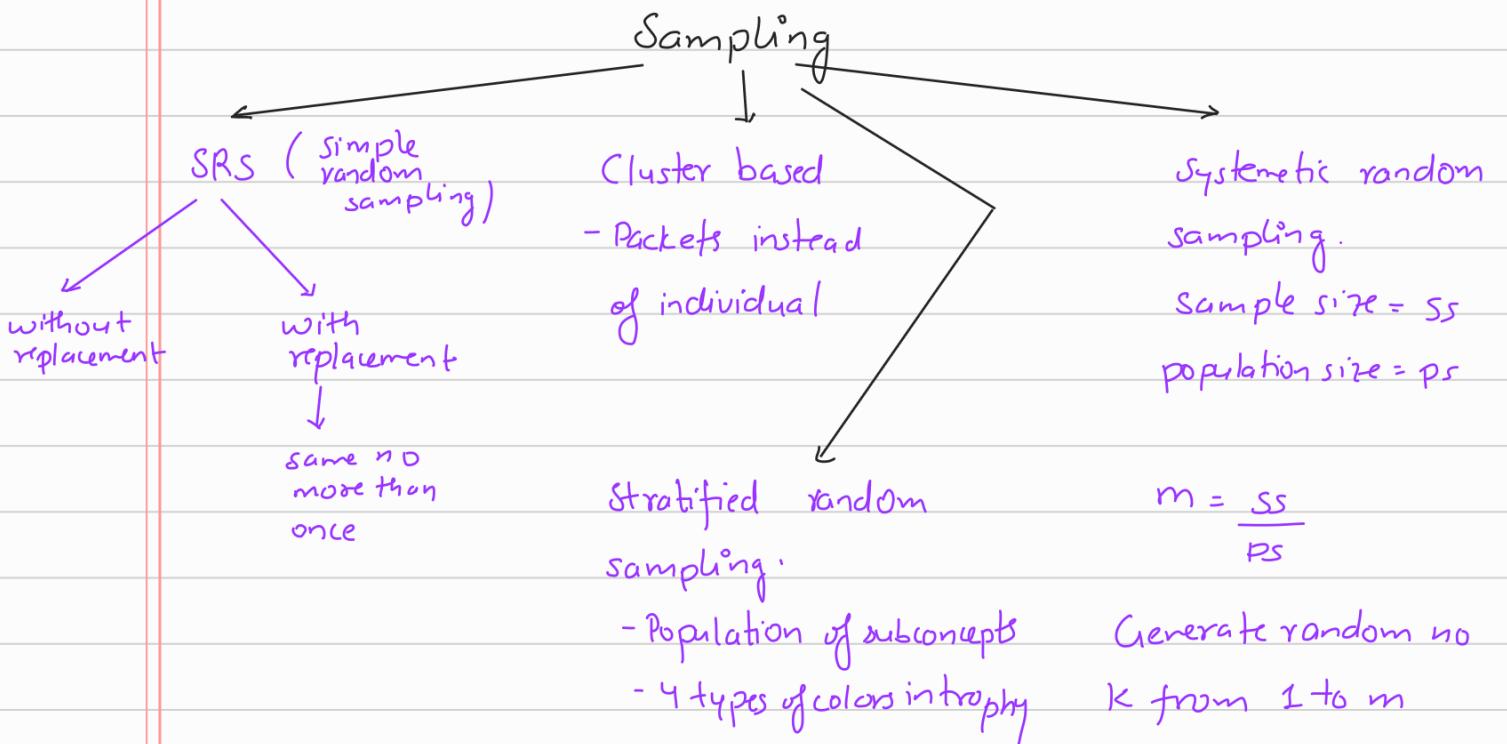
- Q. Driving distance of golf ball - 4 Brands golf balls
 there are 40 golfers 8 are women & 32 are men

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
 $k, k+m, k+2m, \dots$

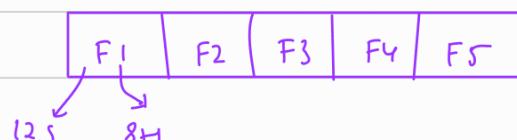
- Cross validation used stratified random sampling
-



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

$c_1 \quad n_1$

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

$c_2 \quad n_2$

$c_3 \quad n_3$

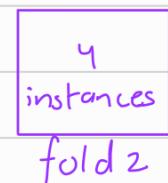
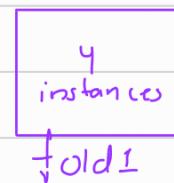
$\vdots \quad \vdots$

$c_m \quad n_m$

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

2-fold cross validation

↪ divide in 2 equal parts



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 z'

S y'

Test

P T n

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

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

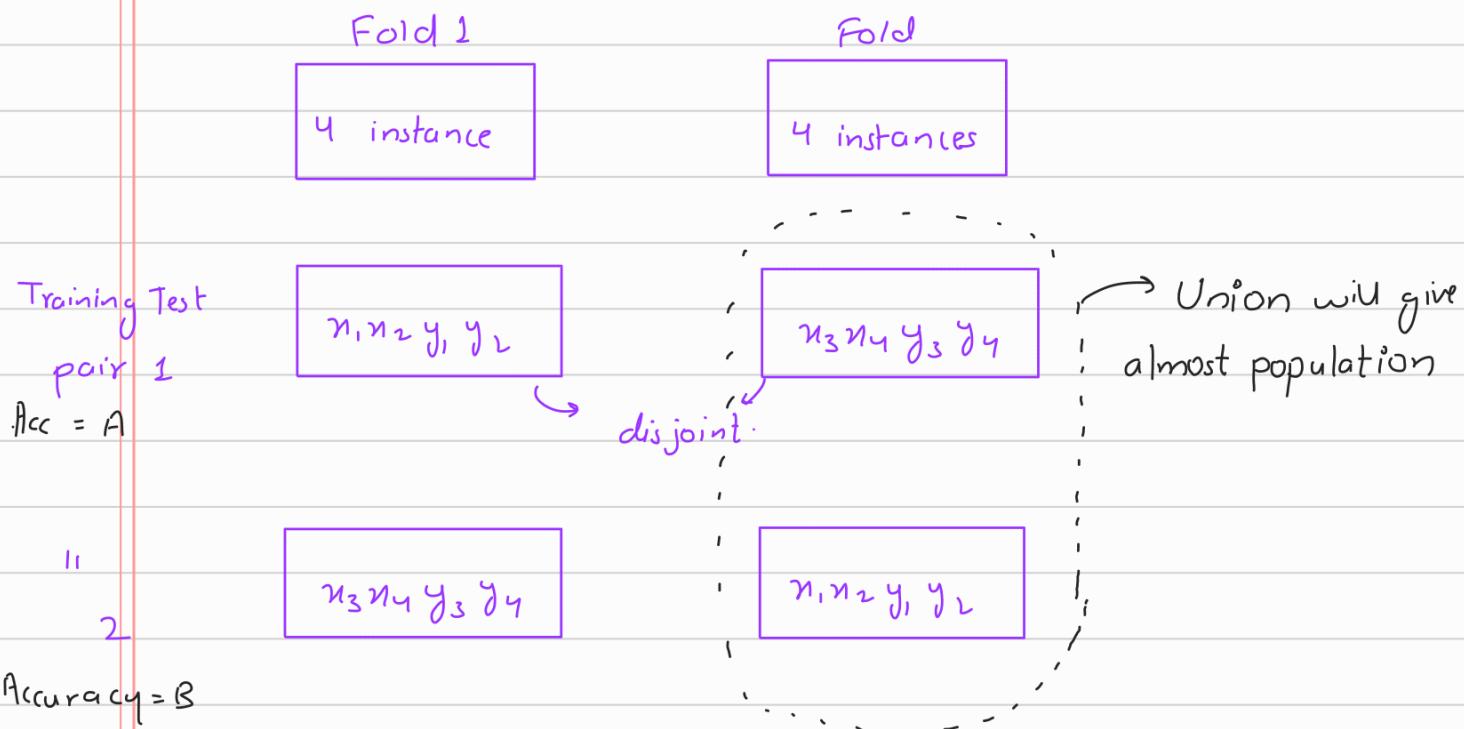
T T $t'_{3''}$

T T $z'_{7''}$

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

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

- 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 2k-fold 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)

Median

(Parent with kid in swimming class)

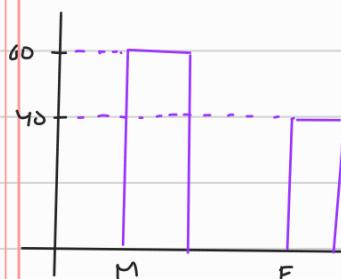
Mode

(5 5 36 36)

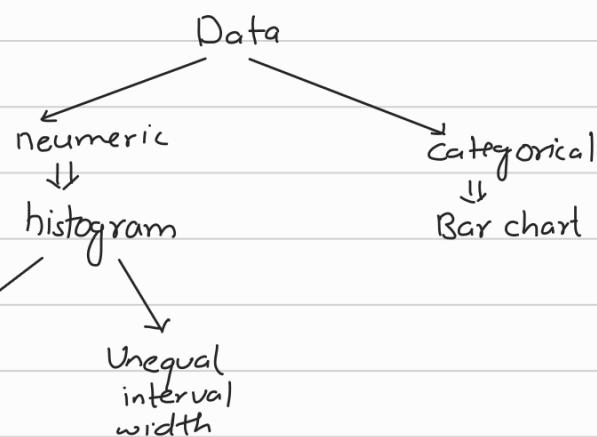
$$\text{mode} = \frac{36+5}{2} = 20.5$$

Distribution →

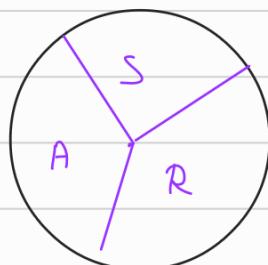
- Bar-chart →



Height of bar \propto Freq.

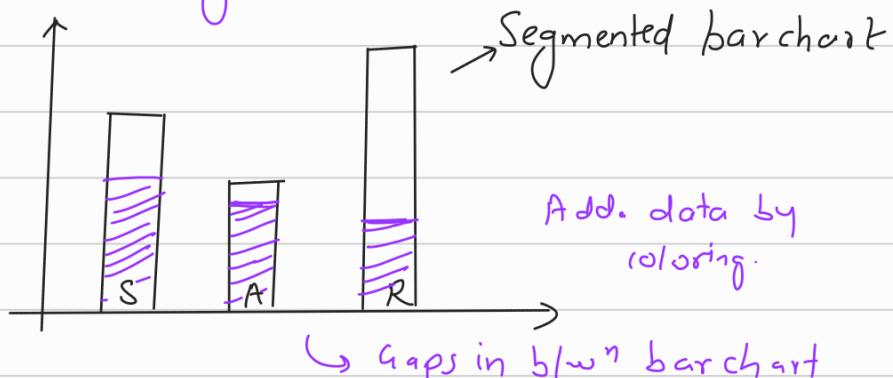


- Pie-chart →



How to plot additional data

- Using barcharts



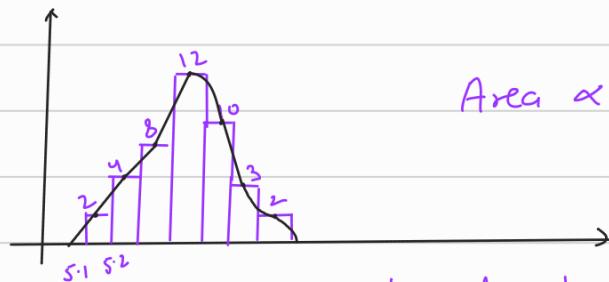
- Histogram → freq

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/wⁿ 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

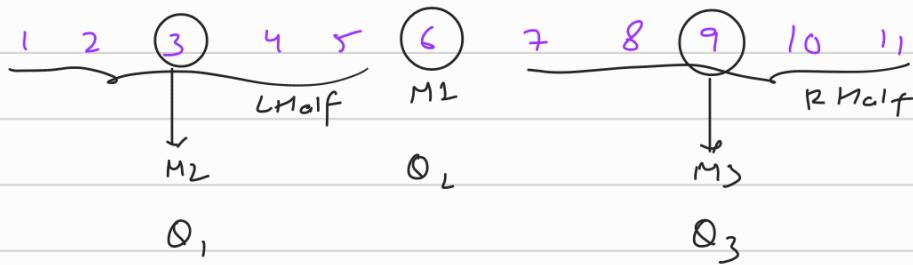
- Mean (fail in outlier)
- Median (fail ^{kid + parent} major class)
- Mode

Variance

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

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

IQR →



25% data lie below Q₁

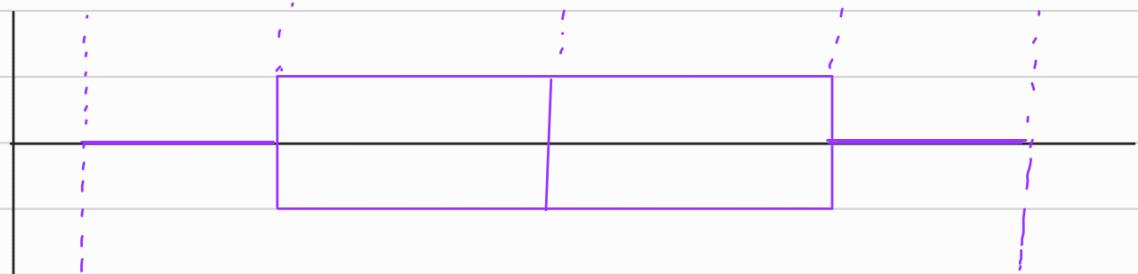
lower bound = 1

Upper bound = 11

Q₁, Q₂, Q₃ = 3, 6, 9

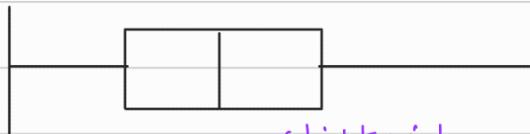
Median = 6

1 2 3 4 5 6 7 8 9 10 11



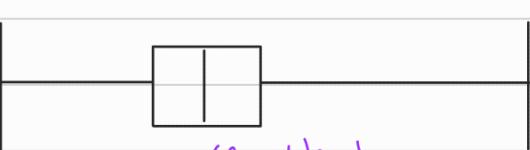
Q.

A



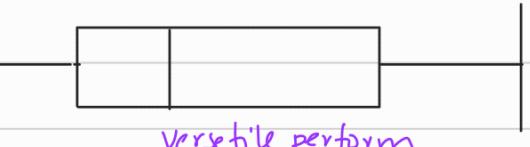
slight risky

B



Whom to select ??

C



versatile perform

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

- What to use

Piechart

Barchart

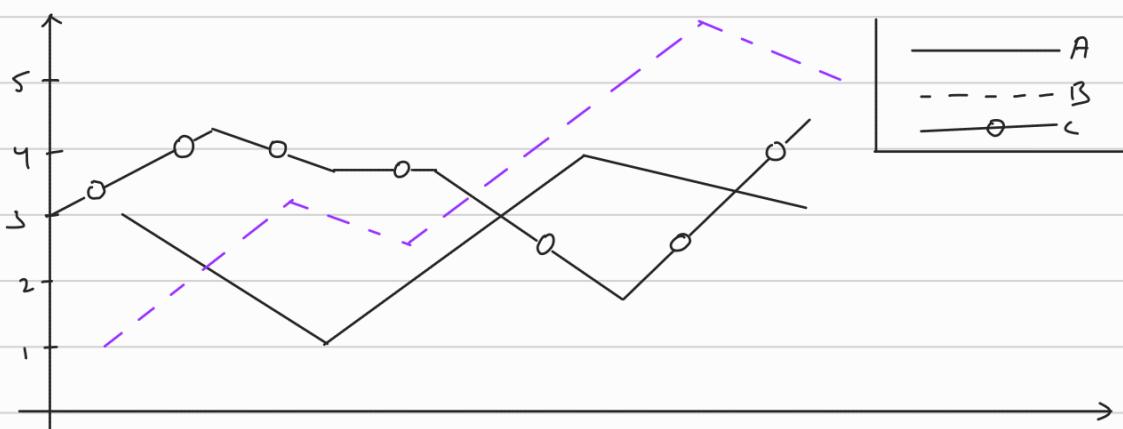
best

line Chart → (diff line similar to below one)

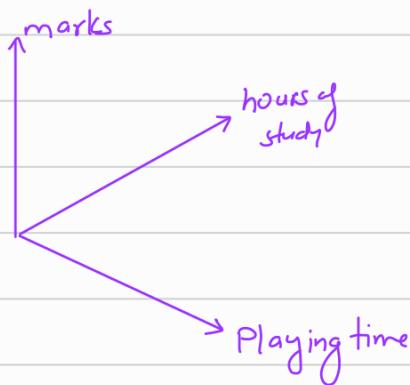
Scatter plots → (diff symbol) for diff people

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 probability 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	summarize	yes (9)	2/9	4/9	3/9
Mild	Mild		No (5)	3/5	2/5	9/5
Mild	hot					
Hot						

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

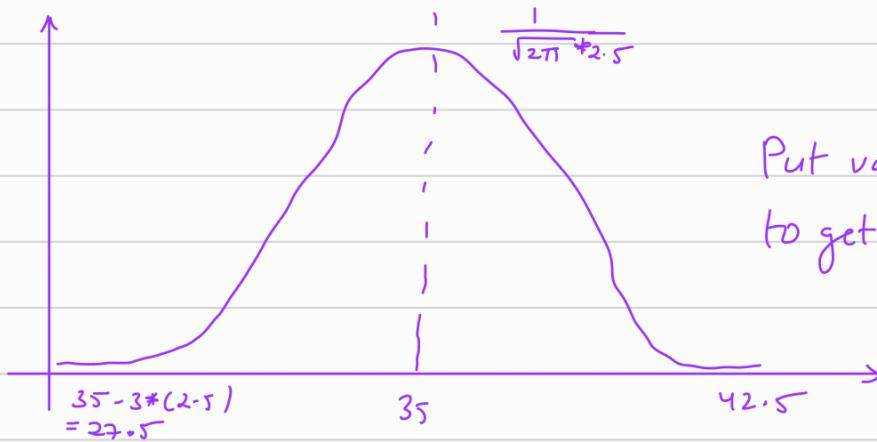
$$f(u) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(u-\mu)^2}{2\sigma^2}}$$

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

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

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

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



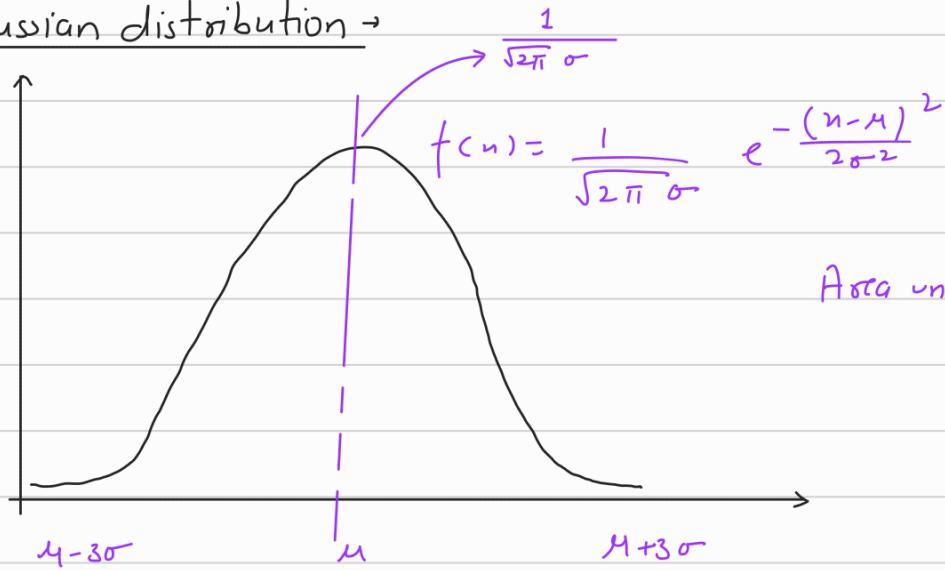
Put value in function
to get ans.

- 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



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 $\mu + 3\sigma$

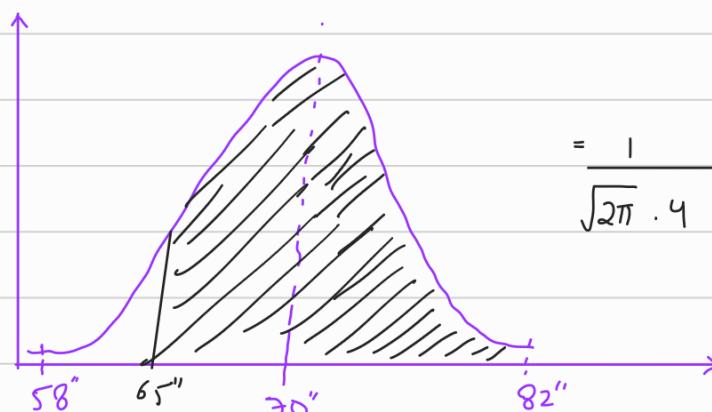
Standard normal distribution \rightarrow (Z-table)

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".



$$= \frac{1}{\sqrt{2\pi} \cdot 4} e^{-\frac{(n-70)^2}{2(16)}}$$

integrate in range

or
transform it to standard normal distribution

$$n = 14$$

$$1 = -4.5$$

to convert σ to 1

$$2 = -3.5$$

divide $n-1/\sigma$

$$3 = -2.5$$

$$4 = -1.5$$

$$5 = 0.5$$

$$6 = -0.5$$

$$7 = 1.5$$

$$8 = 2.5$$

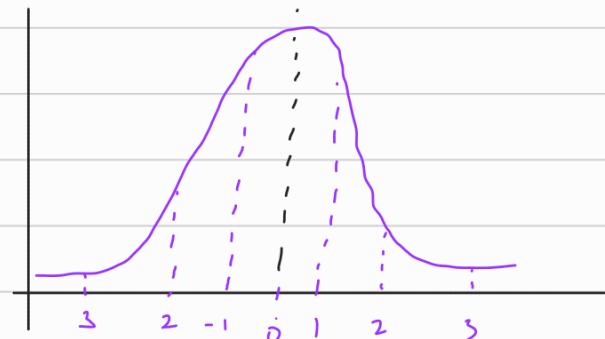
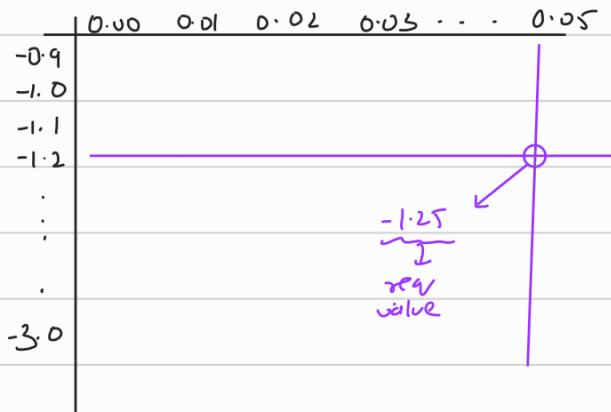
$$9 = 3.5$$

$$10 = 4.5$$

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

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

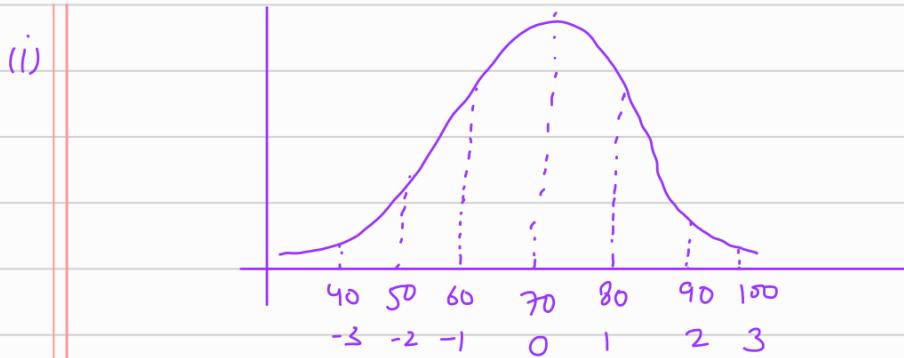
- z-table give area on left hand side of curve



- 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.
 $P(X \geq 80)$, $P(60 \leq X \leq 80)$, $P(X \leq 50)$

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



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

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

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

Binomial distribution →

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

$$P(X=2) = nC_2 \left(\frac{1}{2}\right)^2 \left(\frac{1}{2}\right)^2$$

$n=7$
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

then $E(Y) = aE(x) + b$

$\text{var}(Y) = a^2 V(x)$

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

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

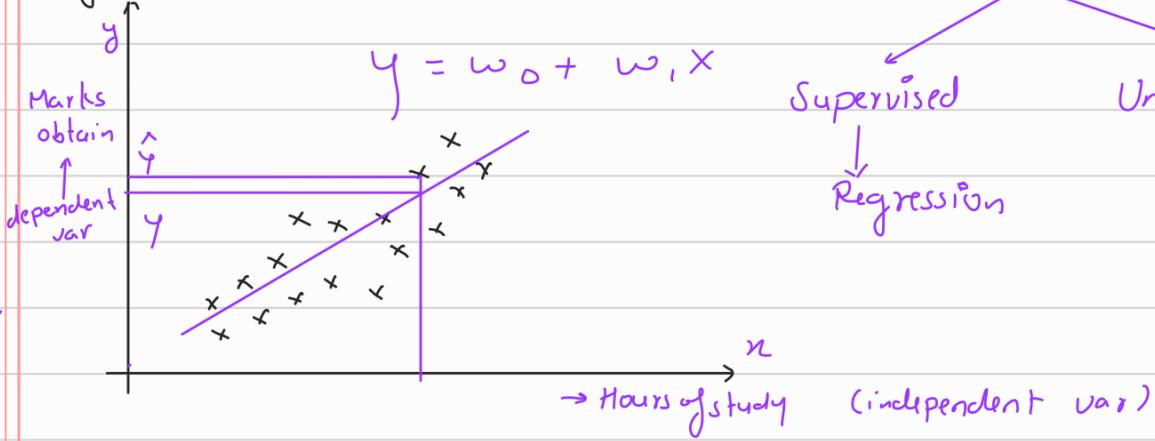
$$= q/p^2$$

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

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

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

Regression →



Machine learning

Supervised

Unsupervised

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 2nd order relationships

$$Y = w_2 \cdot x^2 + w_1 \cdot x + w_0 \quad (\text{reg with 2nd order poly})$$

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	let $y = x$	Residual error	$E = \frac{1}{2} \sum \text{sqv}$	let $y = 2n$	γE	ζw
2	1	1	1	$\xrightarrow{\text{sqv}}$ 1	2	0	0
4	2	2	2	$\xrightarrow{\text{sqv}}$ 4	4	0	0
6	3	3	3	$\xrightarrow{\text{sqv}}$ 9	6	0	0
8	4	4	4	$\xrightarrow{\text{sqv}}$ 16	8	0	0

$E = 15$
this need to min.

$\therefore E=0$

- 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	$11 - 5w$

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

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

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

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

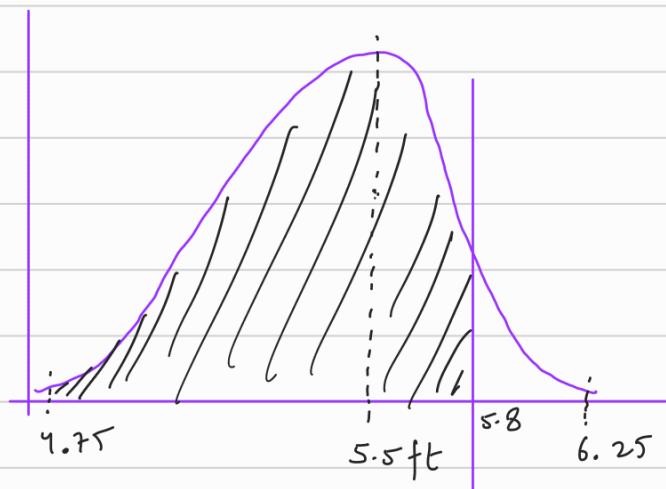
$$= \sum_{i=1}^N y_i n_i - w \sum_{i=1}^N n_i^2 = 0$$

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

$$\Rightarrow \omega = 2$$

$$\therefore y = 2n$$

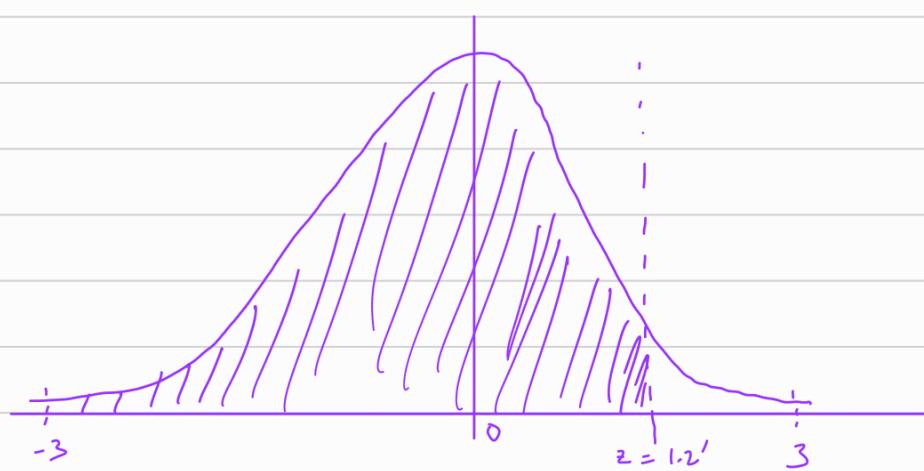
- Q. The height of students of class follow normal distribution with avg height of 5.5ft and variance of $1/16$ ft find out prob that student will have height less than 5.8 ft?



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

$$\sigma = 1/4$$

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

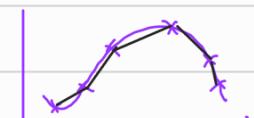


$$z = \frac{5.8' - 5.5'}{0.25} = 1.2'$$

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

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

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

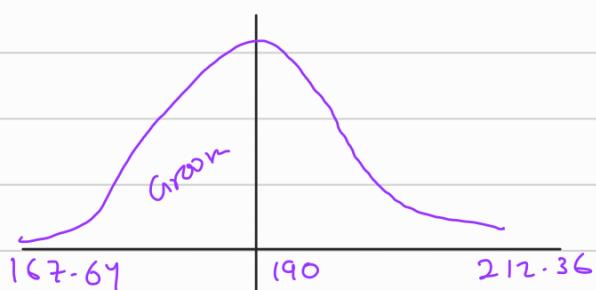
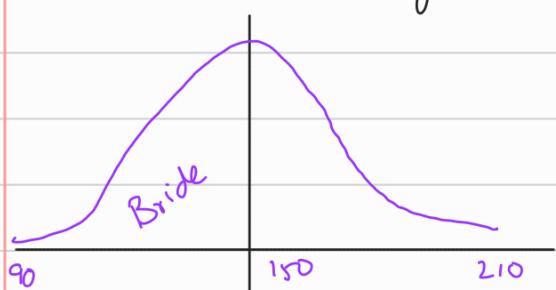
- In regularization we need to tune C & d (check all poss. for best result)
- KEEL DATASET Repo. → regression prob.
- Find 3-4 reg. prob. & use suitable tool box & solve these problems
find optimal value of degree of polynomial & regularization polynomial

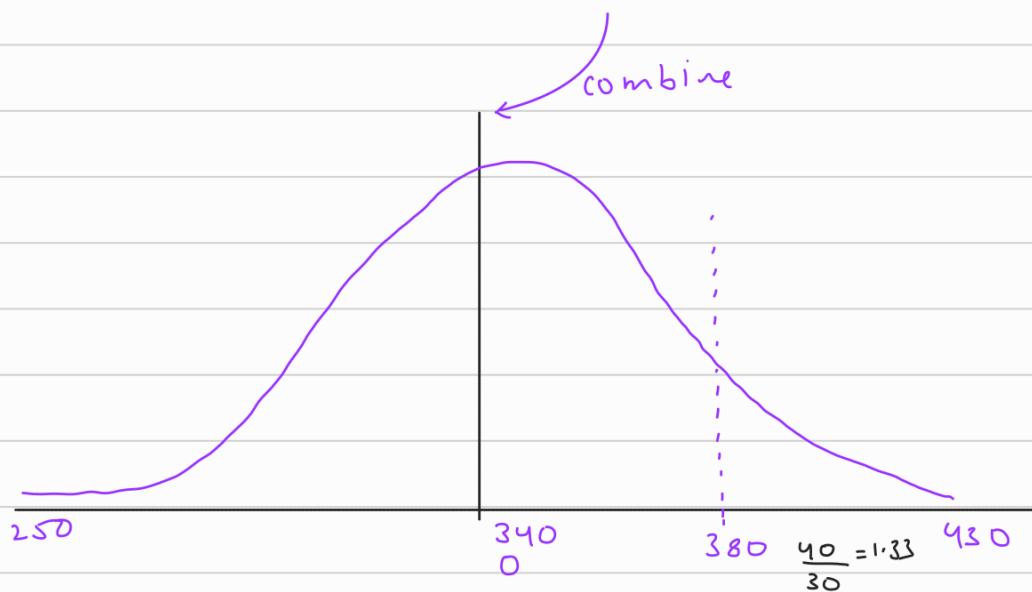
Chap-9 Beyond normal distr. →

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 ≤ 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

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

$$X - Y \geq 5$$

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

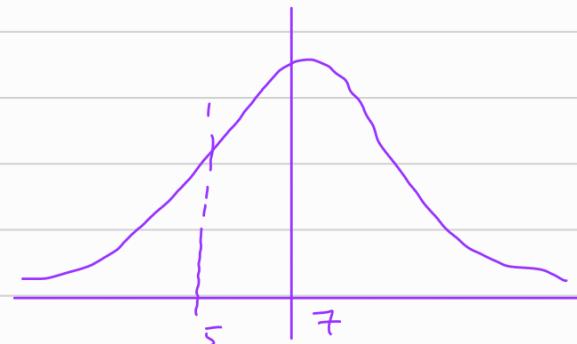
$x, y \in$ independent

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

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

$$E_N = 7, V_N = 36.25$$

$$\sigma_N = 6.02$$



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

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

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 ≤ 30 correct?

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

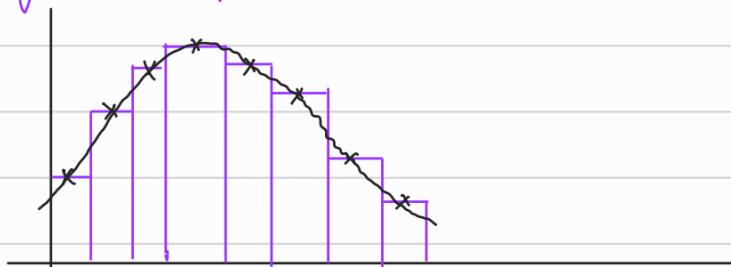
$$P(X \leq 30) = \sum_{r=0}^{30} {}^n C_r p^r q^{n-r}$$

()
→ 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 ' σ '
- 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

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

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

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

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

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

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

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

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

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

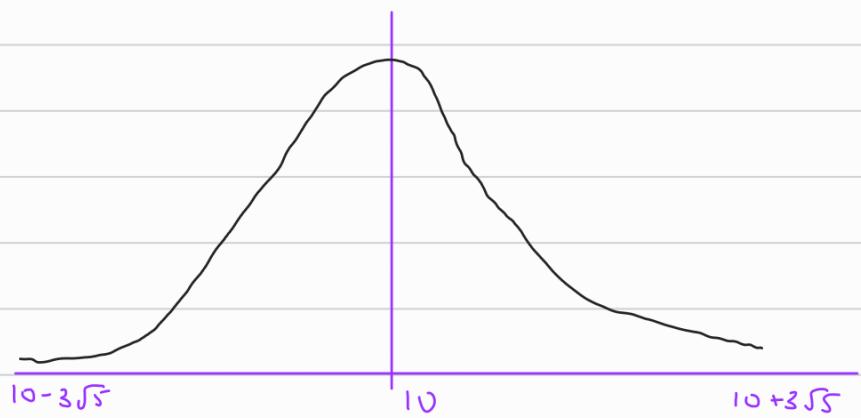
$$= \frac{1}{2^{20}} \cdot (6196) = 0.0059$$

- 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

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

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