

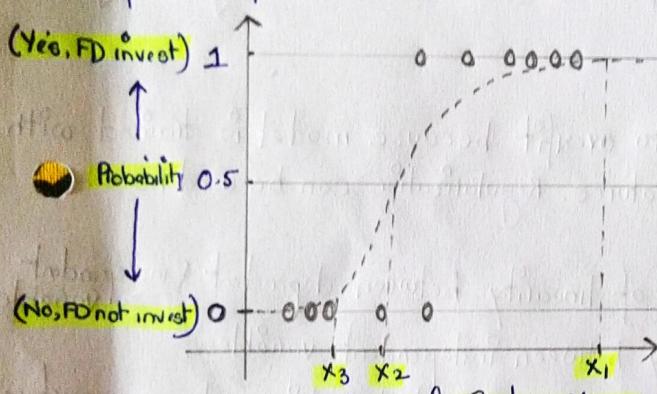
Q) What is logistic regression? Give an example!

Ans - Logistic regression predict the categories (categorical values) instead of continuous values like in Simple / Multiple linear regression.

- LR is a special type of Generalized Linear Model (GLM).

- So consider a dataset where response variable falls into one of the two categories. LR models the probability that Y belong to particular category.

Example → If a customer will invest in Fixed Deposit or Not?



→ LR fits "S" shaped logistic function.

→ The curve goes from 0 to 1, tells us what probability that the customer will invest in fixed deposit.

→ If the customer age is high (x_1), he/she will invest in fixed deposit (high probability).

→ If the customer age is low (x_3), he/she will not invest in FD (Low probability).

→ If the customer age is intermediate (x_2), there is 50% chance will invest in FD.

So, in short if the probability is greater than 0.5/50%, he/she will invest in FD.

Normally, linear regression $\Rightarrow y = b_0 + b_1 * x$; add sigmoid fn on left side.

$$\ln\left(\frac{P}{1-P}\right) = b_0 + b_1 * x, \text{ known as logistic function / logit fn.}$$

$\frac{P}{1-P} = \text{odds}$, that can take any value between 0 and ∞ , then apply log.

$$\log\left(\frac{P}{1-P}\right) \text{ also known as logit function. Also, odds} = \frac{P(\text{event})}{1-P(\text{event})} = \frac{\text{probability event occurs}}{\text{probability event not occurs}}$$

Q) How to create S shape? i.e., Sigmoid function / Logistic function.

Ans - Sigmoid function is a mathematical function having a characteristic that can take any real values and map it between 0 to 1, shaped like letter "S".

- Sigmoid function also known as logistic function.

$$\text{Formula, } y = \frac{1}{1+e^{-z}}$$

- So if value of z goes to positive infinity then predicted value of y will become 1 and if goes to negative infinity then predicted value of y become 0.

- If the outcome of sigmoid function is more than 0.5 then we classify label as class 1 and if it is less than 0.5 then we classify it as class 0.

Q) Give the advantages and disadvantages of Logistic Regression?

Ans - ADVANTAGES OF LR -

- i) It makes no assumption about distribution of target class.
- ii) It measures how appropriate a predictor size (coefficient size) & degree of association.
- iii) It can interpret model coefficient as indicator of feature importance (positive/negative).
- iv) LR is less inclined to overfit but it can overfit in high dimensional dataset.
We may consider (L1 and L2) technique to avoid overfitting.
- v) This algorithm can easily be extended to multi-class classification (more than 2 class target) using softmax classifier, known as Multi-class logistic regression.

DISADVANTAGES OF LR -

- i) On high dimensional data, LR tends to overfit because model is trained with little training data and lots of features. Regularization can be used here.
- ii) LR is sensitive to outliers.
- iii) Major limitation of LR is assumption of linearity between dependent & independent variable.
- iv) LR requires average/no multicollinearity between independent variables.
- v) Non linear problem cannot be solved with LR since it has linear decision surface.
So transformation of non linear feature is required, so data become linearly separable.

Q) What are the assumptions in Logistic Regression?

Ans - Assumptions of LR are -

- i) Outcomes should be categorical in binary or multinomial.
- ii) There should be linear relationship between Independent and dependent variable.
- iii) There should be NO influential value (extreme/outliers) in predictors.
- iv) There should be NO high correlated intercorrelation/multicollinearity among the predictor.

INTERVIEW DAY - ODDS, log(Odds), Odds Ratio, Log(Odds Ratio)

Q) What is Odd? How it is different from probability?

Ans -

$$\text{Odds} = \frac{\text{Something happening}}{\text{Something not happening}}$$

$$\text{Probability} = \frac{\text{Something happening}}{\text{Everything that could happen}}$$

Example 1 - We might say that the odds in favour of my team winning the game are 1 to 4 times.

$$\begin{array}{c} \text{Odds} = \frac{1}{4} = 0.25, \text{ Probability} = \frac{1}{5} = 0.20 \\ \uparrow \quad \nwarrow \uparrow \uparrow \uparrow \\ \text{WIN} \quad \text{LOSE} \end{array}$$

Example 2 - We may say that odds in favour of my team winning the game 5 to 3.

$$\begin{array}{c} \text{Odds} = \frac{5}{3} = 1.7, \text{ Probability} = \frac{5}{8} = 0.625 \\ \backslash \quad / \quad \backslash \quad / \\ \text{Win} \quad \text{Lost} \end{array}$$

Odds are not probability. Odds are ratio of something happening divide by something not happening. Probability is ratio of something happening divide by everything that could happen. In many cases, $\text{odds} = \frac{P}{1-P}$

$$\text{It is derived from, odds(win)} = \frac{\text{Probability(win)}}{\text{Probability(lose)}} = \frac{\text{Probability(win)}}{1-\text{Probability(lose)}} = \frac{P}{1-P}$$

Q) What is log(Odds)? Give an example. And why we added logs.

Ans - Log of odds solve the problem of symmetry.

Suppose my team is good,

$$\text{Odd} = \frac{1}{4} = 0.25$$

more worse, $\frac{1}{8} = 0.125$.

more more worse, $\frac{1}{16} = 0.06$

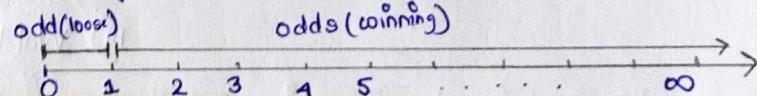
Full worse, $\frac{0}{\text{anything}} = 0$

Suppose my team is good,

$$\text{odds} = \frac{5}{3} = 1.7, \text{ improve } \frac{9}{3} = 3$$

More improve $\frac{27}{3} = 9$, will go to ∞ .

So, range of winning if team is strong is 1 to ∞ .
(Numerator > Denominator)



Range will be 0 to 1 (Denominator > Numerator)

So asymmetry make it difficult to compare odd for ok against my team winning. Example, if odds are against 1 to 6, then $\text{odds} = \frac{1}{6} = 0.17$ but if in favour, $\frac{6}{1} = 6$.

Taking log of this odd can solve and make everything Symmetry.

$$\text{If odds are against 1 to 6} = \log(\text{odds}) = \log(\frac{1}{6}) = \log(0.17) = -1.79.$$

$$\text{If odds are against 6 to 1} = \log_3(\text{odds}) = \log(\frac{6}{1}) = \log(6) = 1.79.$$

this is known as the log of ratio of the probabilities is called logit function, and form basis of logistic regression. Logs make the range $-\infty$ to ∞ , and midpoint is 0.

$$\log(\text{odds}) = \log\left(\frac{P}{1-P}\right)$$

Q) What is Odd Ratio and log(Odd Ratio)? Give an example.

Ans - When people say about odd Ratio, they are talking about ratio of odds.

$$\text{Odd Ratio} = \frac{\text{Odd}_1}{\text{Odd}_2} = \frac{xx}{0000} = \frac{2/1}{3/1} = 0.17$$

So when we calculate odd Ratio of something.

- If denominator is larger than numerator, odd ratio goes from 0 to 1
- If numerator is larger than denominator, odd ratio goes from 1 to ∞ .

Therefore taking, logs(Odd Ratio) will make things symmetrical.

Example, if there is any relationship between Investment in FD (yes/no) and age of customer (high, age ≥ 60 & low, age < 60)

		INVESTED FD	
		Yes	No
Age	High	23	117
	Low	6	210

Total people $\rightarrow 356$.

Total FD invested $\rightarrow 29(23+6)$, Not Invested $\rightarrow 327(117+210)$

Age, High total $\rightarrow 140(23+117)$, Total Low $\rightarrow 216(6+210)$

Given a person of high age, odds of FD invested $= \frac{23}{117}$.

Given a person of low age, odds of FD invested $= \frac{6}{210}$.

Odds ratio $= \frac{23/117}{6/210} = \frac{0.2}{0.03} = 6.88$, so odds ratio tells us that the odds are 6.88 times greater that someone with high age will invest in Fixed deposit.

So odd's ratio / log(Odd's ratio) tells us if there is strong/weak relationship between two variables, like whether or not Age (low/high) increase odds of FD investment.

Larger value means that ~~metadepression~~ age is a good predictor of FD.

Smaller value means that age is not a good predictor of FD investment.

3 ways to determine odds / log(Odd ratio) \rightarrow i) Fisher Exact Test

ii) Chi-square test (to calculate p value)

iii) Wald test (calculate Confidence Interval, pvalue)

Q) Even if Odd formula is ratio, how it is different from Odds Ratio?

Ans - If we take random numbers, that all those odds up to 100 (for example) and use them to calculate odds it will not be normally distributed.

So, if we log(odds) it will become normally distributed.

And how it is different from Odds Ratio is explained above.

INTERVIEW DAY - Probability, Likelihood, Maximum likelihood

Q) What is the difference between probability and likelihood?

Ans - Probability is the percentage that a success occurs. For example, tossing a coin. 0.5 is the probability of success.

- Likelihood is the conditional probability. Same example, we toss coin 10 times and suppose that we got 7 successes and 3 failed.

$$\text{Likelihood}(0.5 | 7) = 0.1171.$$

Meaning, 0.117 is the probability that above event will happen (7 success out of 10 trials) by knowing probability of one success is 0.5 (toss one time).

Therefore,

Likelihood, is the probability (conditional probability) of an event (a set of success) occur by knowing the probability of a success occurs.

Probability, is the percentage that a success occurs.

Another example, in a cricket match, a coin is tossed and one captain calls head and win the toss.

Now, what is the probability that winning captain will elect to bat? $\frac{1}{2}$, either they will elect to bat or bowl, so probability is straight 50%.

But the commentators are discussing what is the likelihood that "winning captain" will elect to bat.

Now that likelihood will not straight to 50% because likelihood depends on type of pitch, strengths and weakness of team, weather etc.

So, probability is straight up mathematics whereas likelihood is a function of many parameters and conditions.

Q) What is Maximum likelihood?

Ans - The goal of Maximum likelihood is to find the optimal way to fit a distribution to the data.

- There are lots of different distributions like normal or exponential.

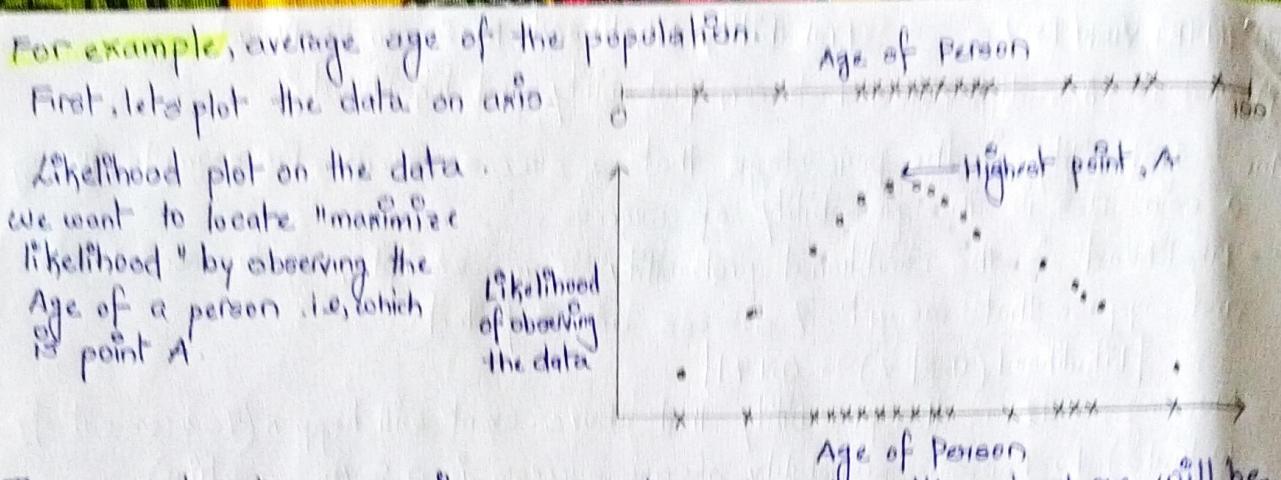
- The reason we want to find / fit a distribution to data is it can be easier to work and more general.

Q) So how Maximum Likelihood works?

Ans - Normally in case of Normal distribution, there are lots of things.

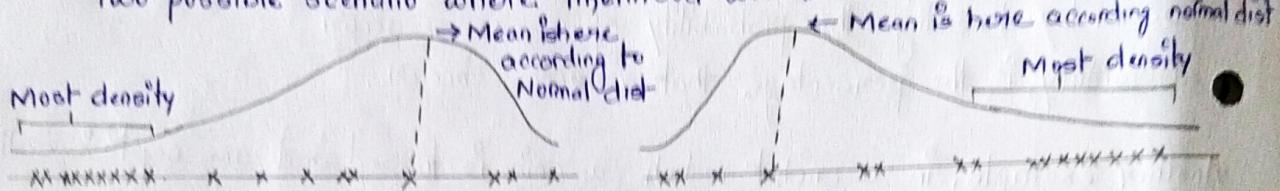
i) We expect most of the measurement to be close to the average/mean.

ii) We expect nearest measurement to be relatively symmetrical around the mean.



There can be other scenario as well, but we have to settle out where will be the center things, i.e., most of the values should be near the average (ideal).

Two possible scenario where likelihood will be low.



Unfortunately, most of the values are measured far from the mean.
So, maximum likelihood find the most optimal way to find distribution of given data.

Q) What is the role of Maximum likelihood in logistic regression?

- Ans. Logistic regression uses Maximum likelihood for parameter estimation.
- In this a probability distribution for target variables (class label) must be assumed and then a likelihood function defines that calculated the probability of observing the outcome given the input data and the model.
 - This function can then be optimized to find the set of parameters that results in the largest sum likelihood over training set.
 - In MLE, we wish to maximize the conditional probability of observing the data (x) given a specific distribution and its parameters.

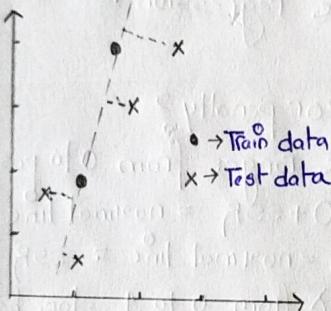
Q) How to address Overfitting in regression?

- Ans - Overfitting, is the model which overfit in training data & performs poorly in test.
- One way could be reduce number of features in the model but there will be loss of information due to discarding. Model will not have benefit of all info.
 - So regularization can be used, regularization is the process which regularize or shrink the coefficient towards zero.

Type of Regularization - i) Lasso Regularization (L1) ii) Ridge regularization (L2)
iii) Elastic Net Regularization.

Q) What is Ridge Regularization (L2 Regularization)?

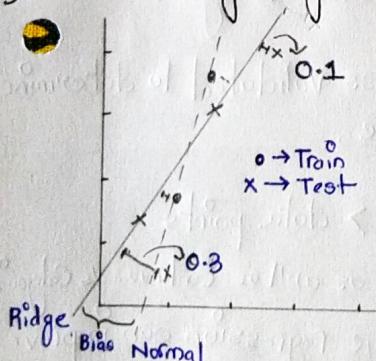
Ans - Suppose we have two datapoints and want to build linear regression.



- Since the line overlaps two data points, minimum sum of square = 0.
- So In train set, SSR = 0 but test set = very high.
- This means line is having high variance, in other word line is overfit in training dataset.
- So Ridge regression comes here, main idea is to find a new line that doesn't fit training data well.

- In other word, we introduce a small amount of bias into line fit to the data.

Q) So how Ridge regression works?



- In normal, it try to reduce/minimize sum of square error.
 - In ridge, it minimize = sum of square residual + λ (slope)².
- Slope², means adds a penalty
 λ , determine how severe the penalty is.

$$y = mx + c, \text{ e.g. size} = y\text{-axis intercept} + \text{slope (weight)}$$

↑
slope

$$\text{Sum of Square (Ridge)} = (0.3)^2 + 0 + 0 + (0.1)^2 = 0.09 + 0.01 = 0.10$$

↑
slope / coefficient

Q) Give an example how normal & Ridge regression?

Ans - Let take example of FD take up and Age.

$$\text{Normal equation, FD} = 0.4 + 1.3(\text{age})$$

here sum of square = 0, slope = 1.3

$$\text{for now } \lambda=1, \text{ penalty} = 0 + 1(1.3)^2 = 0 + 1.69 = 1.69$$

$$\text{Ridge eqn is, FD} = 0.9 + 0.8(\text{age})$$

$$\text{here, SSR} = 0.10, \text{ slope} = 0.8$$

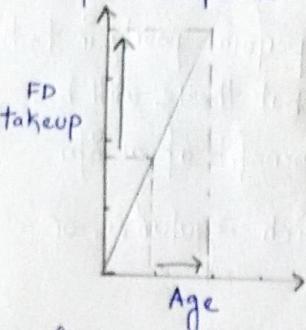
$$\text{for now } \lambda=1, \text{ penalty} = 0.10 + (0.8)^2$$

$$= 0.10 + 0.64 = 0.74$$

So if we want to minimize, sum of square residual plus ridge regression penalty, choose Ridge regression over normal line because Ridge has small bias due to penalty has less variance.

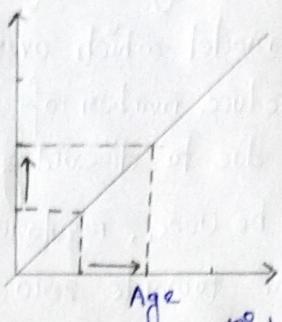
Q) Why Ridge regression is working better than least square method?

Ans - If the slope is steeper



For every one unit increase in age then prediction for FD takeup over two units.

If the slope is normal.



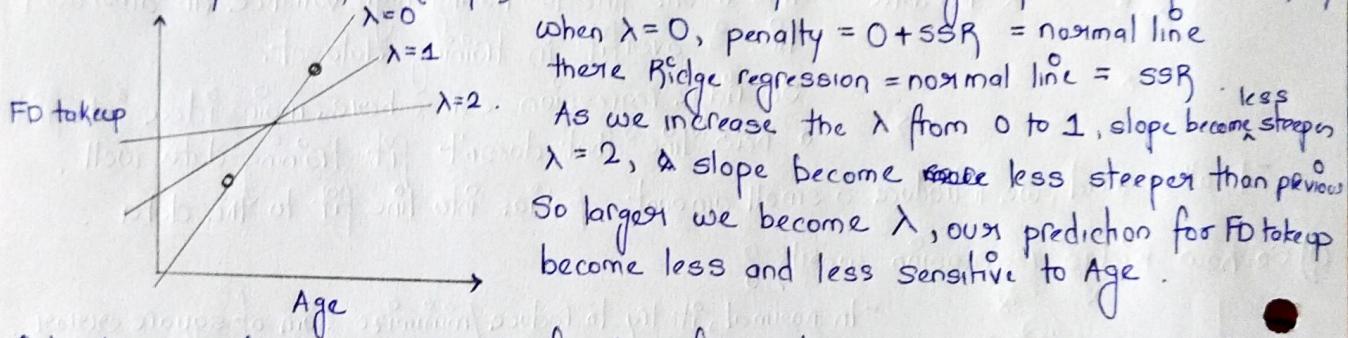
Line suggest that for every one unit increase in Age then there is one unit increase in FD takeup.

So from above, when the slope is ~~steep~~ good / Slope of line is ~~small~~ small, then predictions of FD takeup is less sensitive to change in FD takeup.

Therefore in case of normal vs Ridge, Least Square line is much more steeper than Ridge line. Ridge regression penalty resulted in a line that has a smaller slope which means that predictions made with Ridge regression line are less sensitive to Age than the least square line.

Q) How Ridge regression tries to minimize the error or penalty?

Ans → Penalty = $SSR + \lambda(\text{slope})^2$, λ can be any value ranges from 0 to positive infinity



When $\lambda=0$, penalty = $0 + SSR = \text{normal line}$

there Ridge regression = normal line = SSR

As we increase the λ from 0 to 1, slope becomes less steep

$\lambda=2$, & slope becomes less steep than previous

So larger we become λ , our prediction for FD takeup becomes less and less sensitive to Age.

Optimal value of λ → Use bunch of values of λ and use cross validated to determine which one result in the lowest variance.

Q) So when to use Ridge Regression?

Ans - Ridge regression is used when parameters/features > data points

- Ridge regression works with both Continuous Vs Continuous as well as Continuous Vs Categorical

- In short when sample size are relatively small then Ridge regression can improve predictions (reduce variance) by making prediction less sensitive to training data.

- In general, Ridge regression penalty shrinks but contains all the parameters.

$$\text{Ridge Regression} = \text{sum of squared residuals} + \lambda(\text{slope})^2$$

Ridge regression penalty itself is λ times sum of square parameters except Y intercept & λ is determined using Cross Validation.

Q) What is Lasso Regularization / L1 Regularization?

Ans → Ridge Regression = sum of square residuals + $\lambda (\text{slope})^2$

Lasso Regression = sum of square residuals + $\lambda |\text{slope}|$

Hence also, λ can be any value from 0 to ∞ and determined by cross validation.

Q) So what is the difference between L1 and L2 Regularization?

Ans

L1/Lasso Regularization

L2/Ridge Regularization

- i) Eqn, sum of square residuals + $\lambda (\text{slope})^2$
- ii) Lasso can shrink slope all way to 0.
- iii) Lasso, variable value can be equal to 0

- i) Eqn, sum of square residuals + $\lambda |\text{slope}|$
- ii) Ridge can shrink the slope close to 0.
- iii) Ridge, variable value can be shrunk to very small but cannot be 0.

Example, FD take up = slope + $\beta_1(\text{Age}) + \beta_2(\text{occupation}) + \beta_3(\text{name}) + \beta_4(\text{weight})$

So for FD take up, 2 good variables → Age of a customer, Occupation (Government/Non)

2 bad variables → Name of a customer, Weight of a customer.

So in L1/lasso, value will shrink and useless variable become 0. Only useful variable remain.

In L2/ridge, value will shrink but never $\neq 0$. So all 4 will remain regardless of useless.

Since lasso regression can exclude useless variable from equation, it is little better than Ridge at reducing variance in a model that contain useless variables.

- In contrast, Ridge regression tends to do little better when most variable are useful.

Q) What is Elastic Net Regularization?

Ans - When a dataset include millions of parameters and when we have millions of parameters, we almost need to use some sort of regularization to estimate them. However, variable in those model might be useful/useless, we don't know.

So in those case, we use Elastic Net Regularization which is combination of both Lasso (Some Variable are useful) or Ridge (all variables are useful).

Elastic Net = sum of square residuals + $\lambda_1 |\text{var}_1| + \dots + |\text{var}_n| + \lambda_2 (\text{var}_1)^2 + \dots + (\text{var}_n)^2$

$\lambda_1 \rightarrow$ lasso penalty, $\lambda_2 \rightarrow$ ridge penalty, cross validation for different combination $\lambda_1 \times \lambda_2$ to find

$\lambda_1 = 0$ and $\lambda_2 = 0$, original least square. $\lambda_1 \neq 0$ and $\lambda_2 = 0$, ~~ridge~~ best value lasso

$\lambda_1 \neq 0$ and $\lambda_2 \neq 0$, elastic net $\lambda_1 = 0$ and $\lambda_2 \neq 0$, ridge

- Elastic Net is good when we have highly correlated variables.

$$\text{Elastic Net} = \underbrace{\text{sum of square residuals}}_{\text{Lasso penalty}} + \underbrace{\lambda_1 |\text{var}_1| + \dots + |\text{var}_n|}_{\text{Lasso penalty}} + \underbrace{\lambda_2 (\text{var}_1)^2 + \dots + (\text{var}_n)^2}_{\text{Ridge penalty}}$$

Q) What is Weight of Evidence(WOE)?

Ans - WOE tells the predictive power of an independent variables in relation to the dependent variable. It is generally described as a measure of the separation of good customers and bad customers.

- Good customers refers to the customer who paid loan back and bad customers who defaulted to loan.

$$\text{WOE} = \ln \left[\frac{\text{Distribution of Good}}{\text{Distribution of Bad}} \right]$$

$\ln \rightarrow$ Natural log
 $\text{Distribution of Good} \rightarrow \% \text{ Good customers}$
 in particular group.

$\text{Distribution of Bad} \rightarrow \% \text{ Bad customers}$
 in particular group.

- Positive WOE \rightarrow Distribution of Good $>$ Distribution of Bad.
- Negative WOE \rightarrow Distribution of Good $<$ Distribution of Bad.

In general terms,

$$\text{WOE} = \ln \left(\frac{\% \text{ of non-events}}{\% \text{ of events}} \right)$$

Steps to calculate WOE - i) For continuous variable, split data into bins.

- 2) Calculate the number of event and non events in each bins.
- 3) Calculate the % of events and % non events in each bins.
- 4) Calculate WOE by taking log of division (% non event & % event)

NOTE \rightarrow For categorical variable, we don't need to split the data.

Q) What are the rules related to WOE?

- Ans - i) Each bins should have atleast 5% of observations.
- 2) Each bins should be non-zero for both non-events and events.
- 3) WOE should be distinct for each category. Similar group should be aggregated.
- 4) WOE should be monotonic, either growing or decreasing with grouping.
- 5) Missing values should be binned separately.

Q) Why combine categories with similar WOE?

Ans - Combine categories with similar WOE and then create new categories of an independent variable with continuous WOE values. It is because the categories with similar WOE have almost same proportion of events and non-events. In other word, the behaviour of both the categories is same.

To check correct binning with WOE -

- i) WOE should be monotonic i.e, either growing or decreasing with bins. We can plot WOE and check linearity on graph.
- ii) Perform WOE transformation & check with logistic regression output

Q) What is Information Value (IV)?

Ans - IV is one of technique to select important variable in a predictive model. It help us to rank variable on the basis of importance.

$$\text{Information Value (IV)} = \sum \left(\frac{\% \text{ of non-events}}{\% \text{ of events}} \right) * \text{WOE}$$

INFORMATION VALUE (IV)	VARIABLE PREDICTIVENESS
Less than 0.02	Not useful for prediction (Not useful for modelling)
0.02 to 0.1	Weak predictive power (weak relation, to Good)
0.1 to 0.3	Medium predictive power (medium strength relation)
0.3 to 0.5	Strong predictive power (strong strength relation)
>0.5	Suspicious predictive power (check once again)

- IV is not a good selection method when we have multi class classification (not binary)

- Random forest can detect non linear relationship very well so selecting variable via Information value & using in RF might not produce robust model.

Q) What are the advantages of WOE & IV?

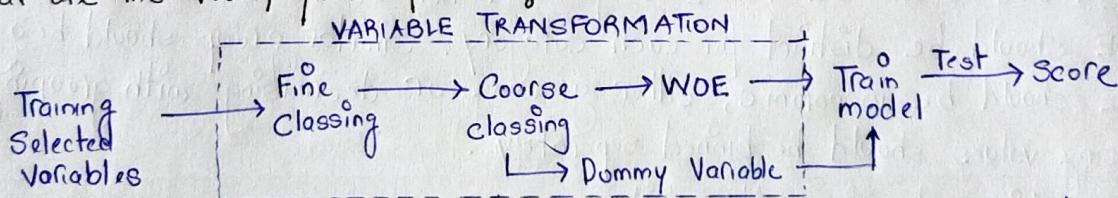
Ans - ① Main practical use of WOE is for encoding, where we can replace the classes with their associated value. For example, suppose in a dataset we found we can replace "Male" with 0.98 and "Female" with -1.58.

② Another positive outcome of WOE is to reduce the number of columns of the input used for training a model. Imagine we have a categorical variable with 10 different classes, one hot encoding will give 10 different columns. Using WOE, classes can replaced by their associated WOE values.

③ As for IV, it provide relationship between independent and dependent variables.

Q) What are the Work flow of WOE?

Ans -



- Fine classing → Applied to all continuous variable and those discrete variable with high cardinality. This is the process of initial binning into typically between 20 and 50 fine granular bins.

To summarize create 10/20 bins for continuous independent variable and calculate WOE & IV of a variable.

- Coarse classing → Combining adjacent categories with similar WOE scores.