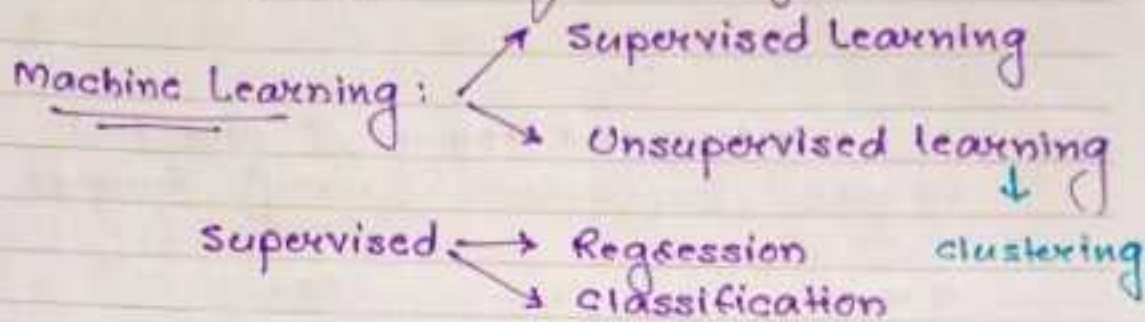


Machine Learning.



Supervised (output is known)

- Regression →
1. Linear Regression
 2. Polynomial Regression
 3. SVR
 4. Decision Tree
 5. Random Forest
 6. Xgboost, 7. KNN

- Classification: →
1. Logistic Regression
 2. SVM
 3. Decision Tree
 4. Random Forest
 5. Naive Bayes
 6. KNN

Unsupervised learning: (output is unknown)

- clustering →
- DBScan, kmeans,
 - Hierarchical,
 - Silhouette scoring

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

Supervised Learning

Eg. Independent features → dependent features
Degree Experience (Salary) features

B.Tech	7	50k
PhD	2	70k

o/p = continuous → Regression problem

Eg. Independent Feat. → dependent
No. of play has No. of study has Pass/Fail Feat.

o/p = categorical = classification problem

- Flight Price Prediction → Regression problem
- Algerian Fire Forest → classification problem
- Predict Air Quality Index → Regression Problem
- Rain Prediction → classification Problem
- Buying day of a person → classification problem

•• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

Note: Independent features is basically input features.

Simple Linear Regression.

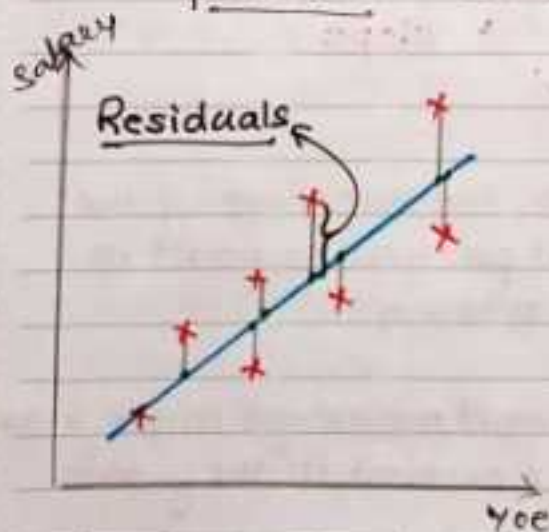
(one independent feature and one dependent feature)

Eg. aim: to create a model, which takes input as height and predict weight.

dataset: height, weight.

Eg. aim: Based on the no. of rooms, predict price.

Eg. model: year of experience and salary
predict: salary band on 1p year.

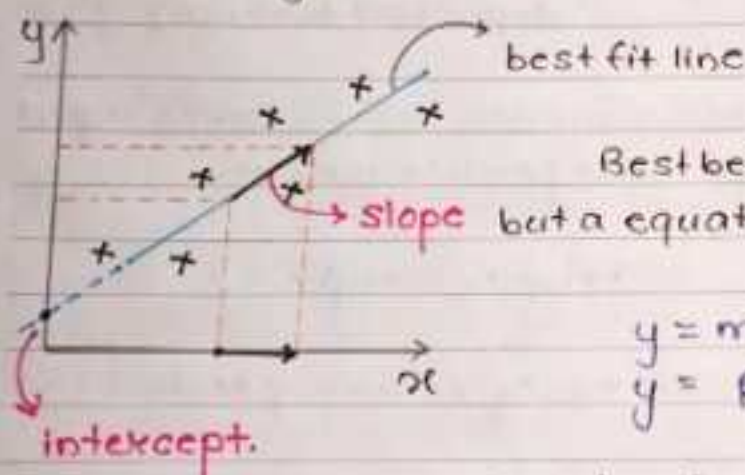


• $\times \rightarrow$ datapoint (original)
• $\bullet \rightarrow$ predicted datapoint.

• difference between real points and predicted points is called residuals or Error.

- Based on the training dataset, it finds the best fit line in such a way that the sum of difference between real points and predicted should be minimum.

First, we need to understand why are creating a straight line?



Best best-fit line is nothing but a equation of straight line

$$y = mx + c$$

$$y = \beta_0 + \beta_1 x$$

$$h_0(x) = \theta_0 + \theta_1 x$$

all are
eqⁿ. of
line.

$$h_0(x) = \theta_0 + \theta_1 x$$

intercept

slope

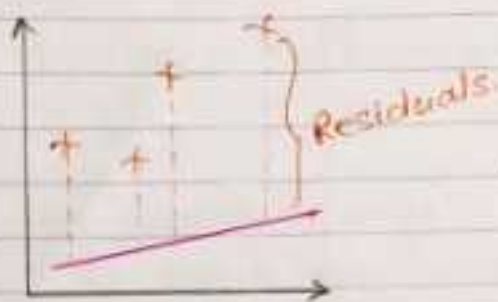
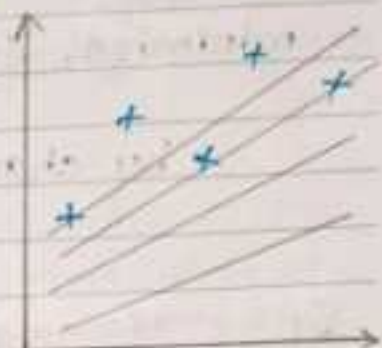
intercept: when $x=0$, the line meeting the y-axis, that particular point is known as intercept.

θ_1 : Slope: with the unit movement in the x-axis, what is the movement in the y-axis.

- By changing θ_0 and θ_1 , best fit line will be change.

θ_0, θ_1 : changing the values, will change the line.

after some iterations, we will get a best fit line which is known as training of the model.



we need to minimize this errors using equation called as **cost function**.

Cost Function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \left(\underbrace{h_{\theta}(x^{(i)})}_{\text{predicted}} - \underbrace{y^{(i)}}_{\text{actual}} \right)^2$$

for easy calculations (derivatives)

mean square error. (one of the cost function)

- we need to minimize cost function to get the best fit line.
- dividing by m to get the average (mean).

Final aim:

$$\text{minimize, } J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h(\alpha)^{(i)} - y^{(i)})^2$$

• by changing the values of θ_0 and θ_1 .

Eg. Training dataset.

x	y
1	1
2	2
3	3

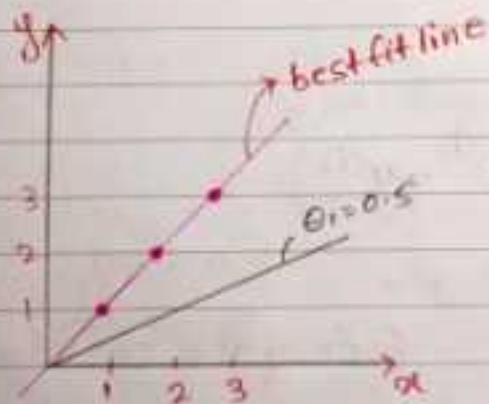
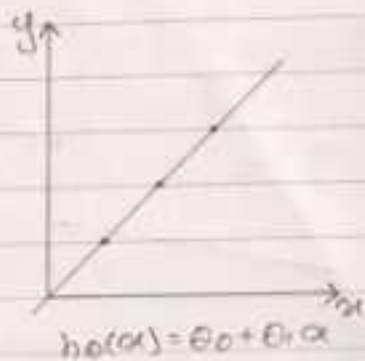
Let us consider,

$$\theta_0 = 0$$

$$h(\alpha) = \theta_1 \alpha$$

Now, let's assume $\theta_1 = 1$

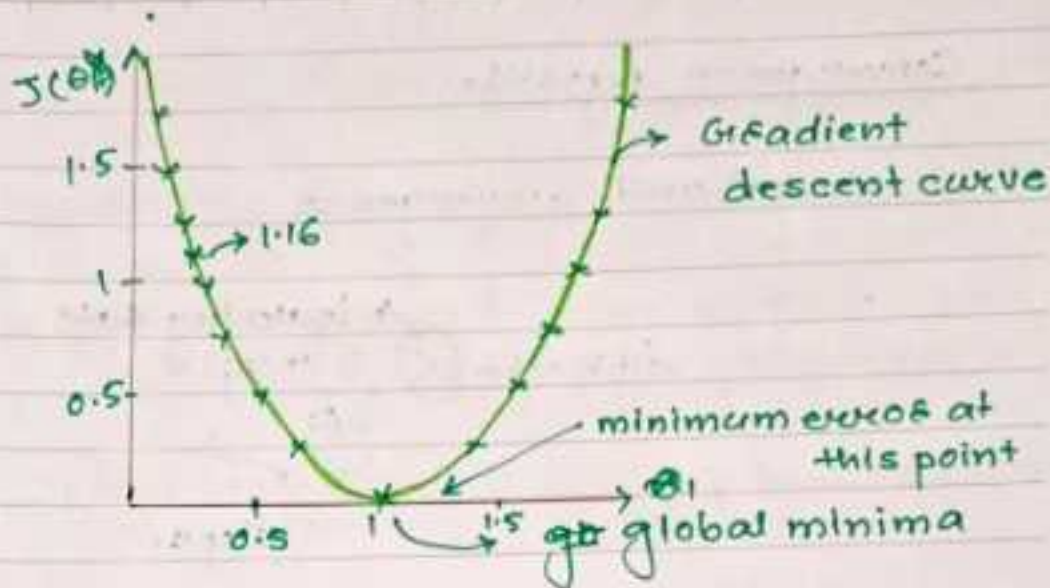
$$\therefore h(\alpha) = \alpha$$



$$\begin{aligned} J(\theta_1) &= \frac{1}{m} \sum_{i=1}^m (h(\alpha)^{(i)} - y^{(i)})^2 \\ &= \frac{1}{3} [0 + 0 + 0] \\ &= 0. \end{aligned}$$

when $\theta_1 = 0.5$

$$\begin{aligned} J(\theta_1) &= \frac{1}{3} [(0.5 - 1)^2 + (1 - 2)^2 + (1.5 - 3)^2] \\ &= \underline{1.16}. \end{aligned}$$



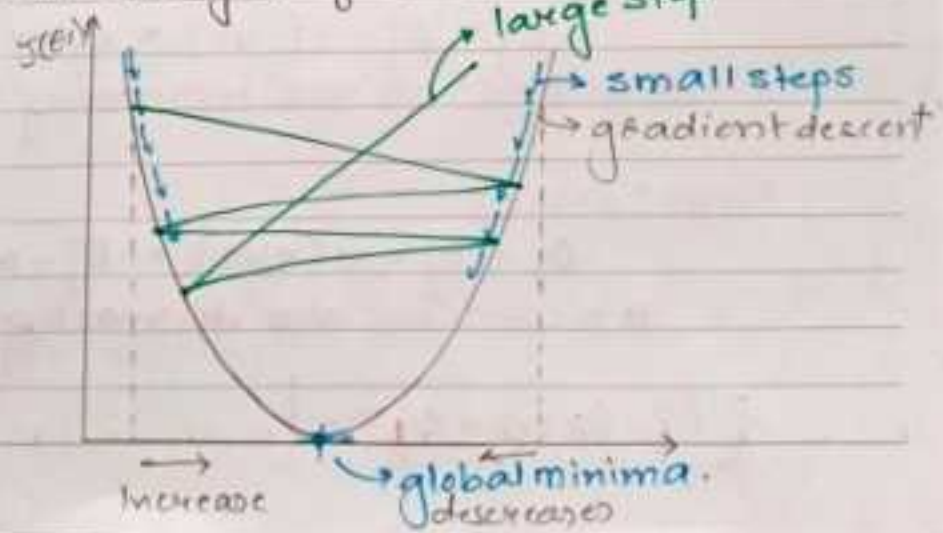
Objective: our main aim is to come near global minima.

we cannot change θ value manually, there should be some mechanism to change θ value to get the global minima.

To overcome this, we use Convergence algorithm.

Convergence algorithm:

- optimize the changes of θ_1 value.



Convergence algorithm:

Repeat until convergence,

{

$$\theta_j = \theta_j - \underbrace{\alpha}_{\text{learning rate}} \underbrace{\frac{\partial J(\theta_j)}{\partial \theta_j}}_{\text{slope}}$$

}

+ve or -ve slope:

Right side of the line facing downwards
→ '-ve' slope

Right side of the line facing upwards
→ '+ve' slope

'-ve' slope:

$$\theta_j = \theta_j - \alpha(-ve) = \theta_j + \alpha$$

It means we are increasing θ_j .

'+ve' slope:

$$\theta_j = \theta_j - \alpha(+ve) = \theta_j - \alpha$$

It means we are decreasing θ_j .

$$\boxed{\theta_j = \theta_0 \text{ and } \theta_j \neq \theta_1}$$

α , Learning Rate

It decides the speed of the convergence.

If α is very small: It will take more time to reach global minima.

If α is very large: It will jump here and there, and won't reach to global minima.

α , should be around 0.001 for smaller steps.

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

convergence algorithm:

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \left(\underbrace{h_{\theta}(x^{(i)})}_{\text{predicted}} - \underbrace{y^{(i)}}_{\text{actual or truth point}} \right)^2$$

mean square error

cost function

m = no. of datapoint.

difference between cost function and loss function.

cost function: we find errors for all the points and take average of it.

loss function: we find errors for observed points and if we find errors for all the points then it is loss function.

$$\begin{aligned} \text{loss function} &= \left(\underbrace{h_{\theta}(x^{(i)})}_{\text{predicted value}} - \underbrace{(y^{(i)})}_{\text{actual value}} \right)^2 \\ &= (\hat{y}_j - y_j)^2 \end{aligned}$$

* for every single point we need to find loss function.

** Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

To achieve global minima,

derivative w.r.t. $\theta_0, j=0$

$$\frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_0} \left\{ \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right\}$$

$$= \frac{1}{2m} \quad h_{\theta}(x) = \theta_0 + \theta_1 x$$

$$= \frac{\partial}{\partial \theta_0} \left\{ \frac{1}{2m} \sum_{i=1}^m [(\theta_0 + \theta_1 x^i) - y^i]^2 \right\}$$

$$= \frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^i) - y^i) \times 1$$

$$\frac{\partial}{\partial x} (1+x-y)^2$$

$$= \frac{\partial}{\partial x} [(1+x-y) \times 1] = \frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^i) - y^i)$$

derivative w.r.t. $\theta_1, j=1$

$$\frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) = \frac{\partial}{\partial \theta_1} \left\{ \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2 \right\}$$

$$= \frac{\partial}{\partial \theta_1} \left\{ \frac{1}{2m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^i) - y^i)^2 \right\}$$

$$= \frac{1}{m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^i) - y^i) x^i$$

Repeat untill convergence

{

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \cdot \sum_{i=1}^m (h(\alpha)^i - y^i)^2$$

$$\theta_1 := \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h(\alpha)^i - y^i) x^{(i)}$$

}

* learning rate: speed of convergence

Types of cost function:

- ① MSE: mean square error
- ② MAE: mean absolute error
- ③ RMSE: root mean square error.

Mean Square Error

$$MSE = \frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2$$

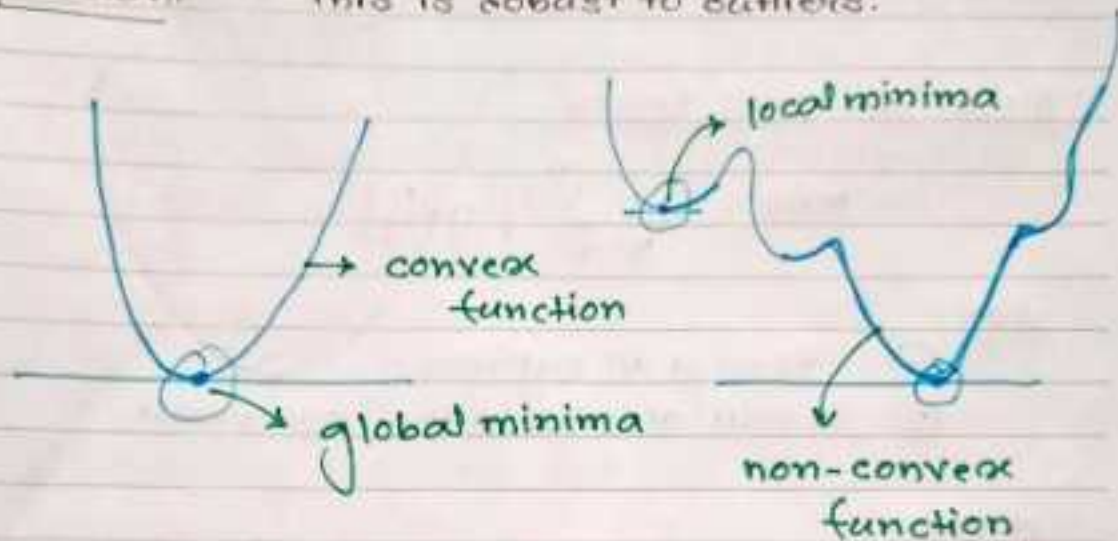
actual.

$$\hat{y} = \theta_0 + \theta_1 x$$

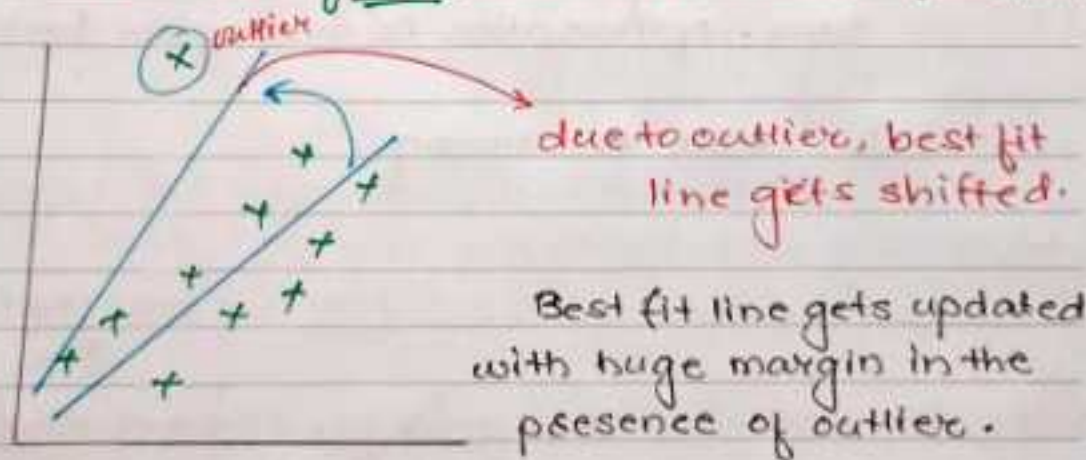
predicted.

- adv
- This equation is differentiable.
 - It also has one global minima.

- disadv
- This is robust to outliers.



∴ main adv. of MSE that it has convex function.



Eg. (salary) → dependent feature
Experience → Independent feature

$$(y - \hat{y})^2 \text{ (lakhs)}^2$$

↪ unit changing → time complexities increases.

we don't do scaling for dependent features.

$(y - \hat{y})^2 \rightarrow$ squared \rightarrow error \rightarrow penalized
 \downarrow
Increased.

Mean Absolute Error

$$MAE = \frac{1}{n} \sum_{i=1}^n |y - \hat{y}_i|$$

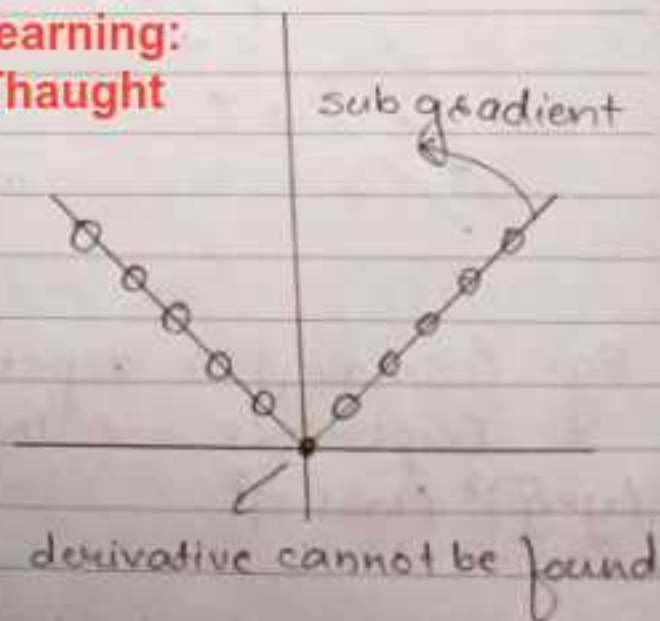
adv:

- ① Robust to outliers.
- ② It will also be in the same unit.

disadv:

- ① convergence usually takes more time. optimization is a complex task.
- ② time consuming.

• Download Machine Learning:
<https://t.me/AIMLDeepThought>



Root mean Square Error (RMSE)

$$RMSE = \sqrt{\frac{1}{n} \sum_{i=1}^n (y - \hat{y})^2}$$

adv:

- It is differentiable.
- unit remain same.

disadv:

- Not robust to outliers.

Huber Loss Function

The Huber loss offers the best of both worlds by balancing the MSE and MAE together.

$$L_{\delta}(y, f(x)) = \begin{cases} \frac{1}{2} (y - f(x))^2 & \text{for } |y - f(x)| \leq \delta, \\ \delta |y - f(x)| - \frac{1}{2} \delta^2 & \text{otherwise} \end{cases}$$

It says: for loss values less than δ , use the MSE, for loss values greater than δ , use the MAE.

using the MAE for larger loss values mitigates the weight that we put on outliers so that we still get a well-sounded model.

At the same time, we use the MSE for the smaller loss values to maintain a quadratic function near the centre.

This has the effect of magnifying the loss values as long as they are greater than 1.

once the loss for those data points dips below 1, the quadratic function down-weights them to focus the training on the higher-error data points.

Note: use the Huber loss any time you feel that you need a balance between giving outliers some weight, but not too much.

How can we check if a model is good or not?
→ using performance metrics.

Performance metrics

- ① R-squared
- ② Adjusted-R squared

R-squared: measures the performance of the model.

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

$$R\text{-squared} = 1 - \frac{SS_{\text{res}}}{SS_{\text{Total}}}$$

SS_{res} = Sum of Square Residuals

SS_{Total} = Sum of square average



$$R\text{-squared} = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

low value

high value

\bar{y} = average of y .

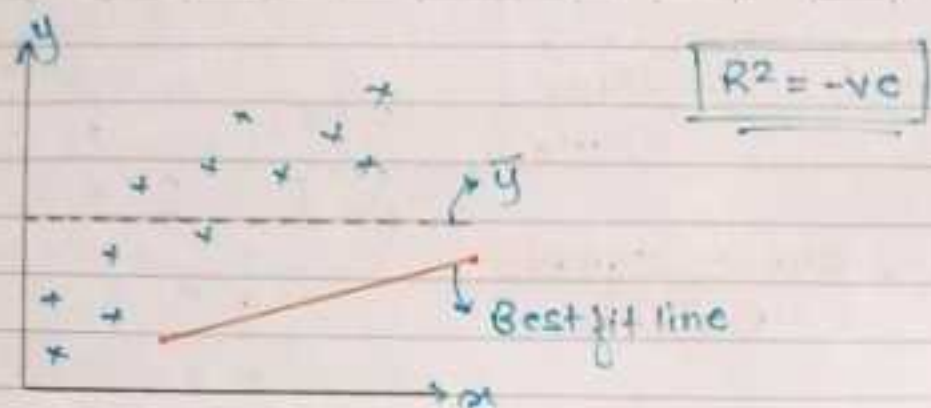
$$= 1 - \left\{ \frac{\text{small value}}{\text{Bigger value}} \right\} \rightarrow \text{small value}$$

$$R\text{-squared} \leq 1$$

If $R\text{-squared} = 0.85 = 85\%$ accurate
 $= 0.75 = 75\%$ accurate

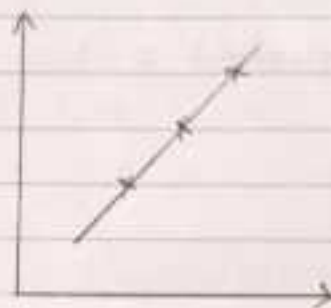
* If $R\text{-squared}$ is 'low' then the model is not working is good.

• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>



Here $(y_i - \hat{y}) > (y_i - \bar{y}) \rightarrow R^2$ will be negative.

$R^2 = 1 \rightarrow$



Adjusted R-squared

Eg.

No direct correlation

size of House	city location	no. of bedrooms	salary	Gender
				Price

Before, when only size of house was present as an independent feature to get the price and has some R-squared value. But with addition of city location R^2 got increase. after that with addition of no. of bedroom, gender the value got increase, but there is ~~not~~ no direct correlation between gender and price, the value shouldn't increased.

R^2	features added	adj. R^2	Independent features
65%	Size of house	63%	$P=1$
75%	location	73%	$P=2$
88%	no. of bedrooms	86%	$P=3$
90%	Gender	85%	$P=4$

no direct correlation with price
 * very slightly increase, but this shouldn't happen.
 To solve this problem we use Adjusted R^2 .

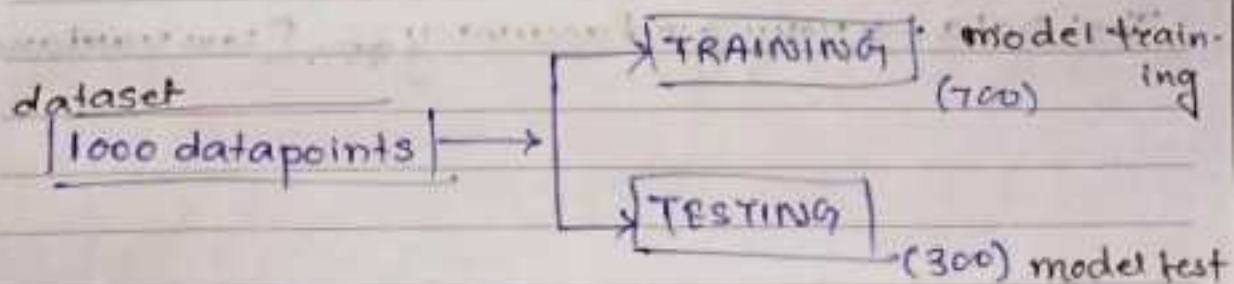
$$\text{Adjusted } R^2 = 1 - \frac{(1 - R^2)(N - 1)}{N - P - 1}$$

N = No. of data points

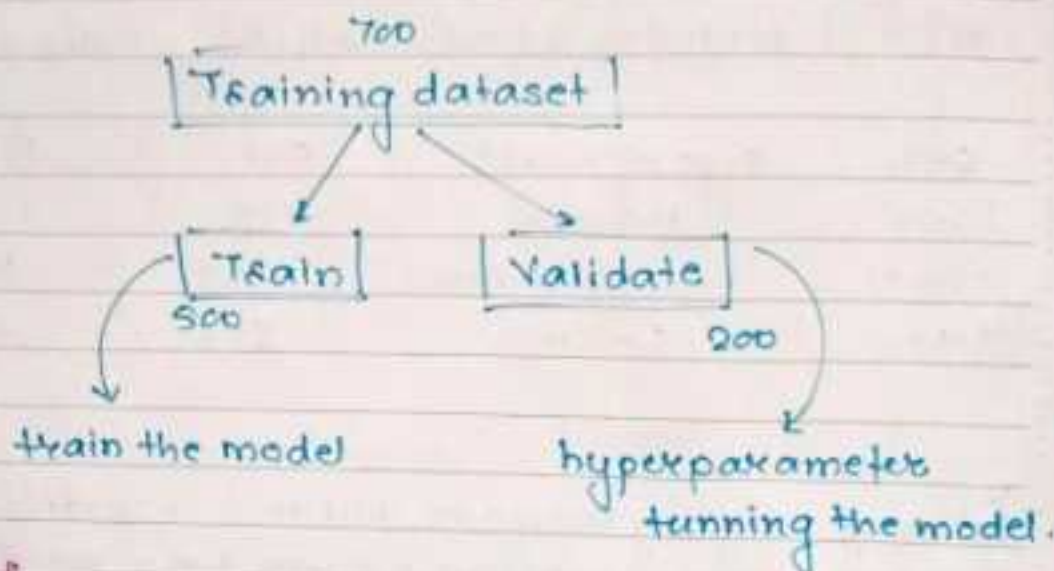
P = No. of independent features

* Adjusted R^2 is the best metrics to evaluate the model.

Overfitting and Underfitting



overfitting and underfitting (Bias and Variance)



MODEL

Generalize Model

Train Data \rightarrow very good accuracy (90%)

Test Data \rightarrow very good accuracy (85%)

our aim is to get good test and train accuracy.

Train Data \rightarrow very good accuracy 90% [low bias]

Test Data \rightarrow very good accuracy 85% [low variance]

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

Overfitting

TRAIN \rightarrow Very good accuracy (sor.) [low Bias]

TEST \rightarrow Bad accuracy (sor.) [High Variance]

Underfitting

TRAIN \rightarrow model accuracy is low. [High bias]

TEST \rightarrow model accuracy is low/high. [low or high variance]

- * we can solve this issue by performing hyperparameter tuning or by increasing no. of dataset.

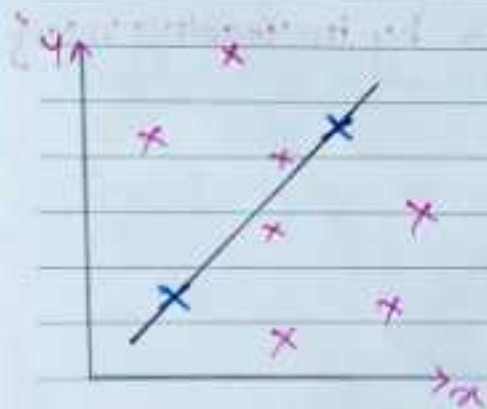
• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

Regularization

Sometimes when we train a model it will start to overfit. A way to avoid overfitting data (especially for models like linear regressions that are heavily affected by outliers) we can use **Regularization**. This will lead to a more general model that is technically less accurate but generalizes to the data better.

1. Ridge / L2 Regression

(used to reduce overfitting)



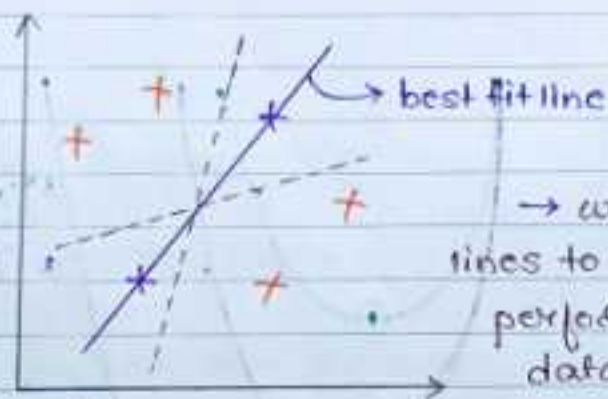
cost function = 0

Training data: low bias

Testing data: low/high variance

- If the data (test data) is near to best fit line then performance will be good.
(low variance)
- If the test data is far (away) to best fit line then performance will be bad.
(high variance)

aim: To reduce overfitting



cost function:

$$\text{cost function} = \frac{1}{m} \sum_{i=1}^m \left\{ h(\mathbf{x}^{(i)}) - y^{(i)} \right\}^2 + \lambda (\text{slope})^2$$

λ = hyperparameter

Eg.

$$h(\mathbf{x}) = \theta_0 + \theta_1 x \rightarrow \text{slope} = \theta_1$$

if multiple features are present, then

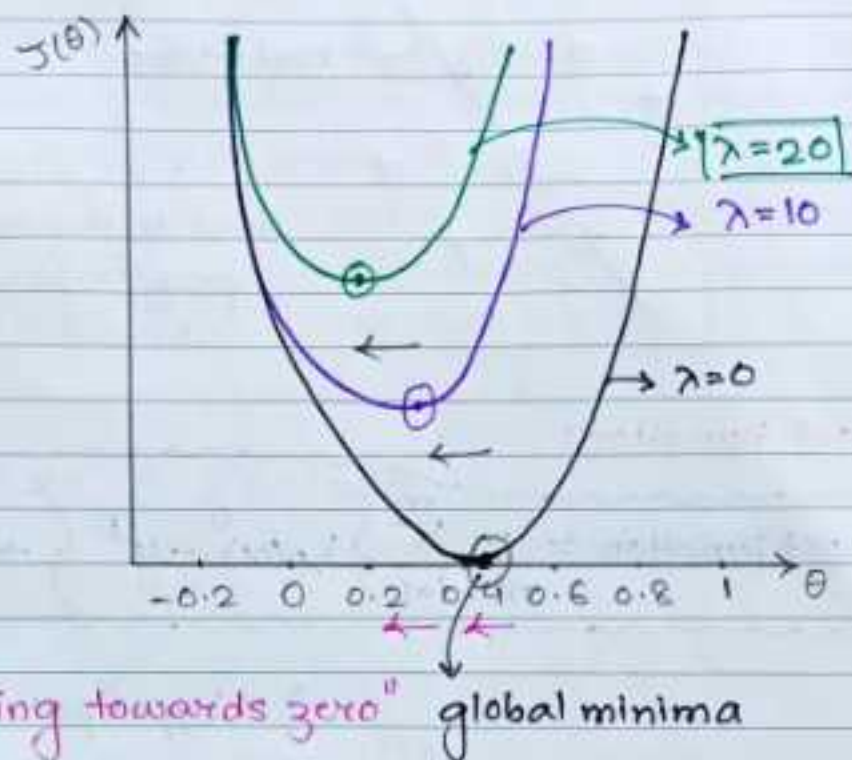
$$(\text{slope})^2 = \sum_{i=1}^n (\text{slope})^2$$

slope = slope of different lines.

if $\lambda = 0$,

cost function is same as linear regression cost function.

Relationship between slope and λ



"shifting towards zero" global minima

- Global minima gets shifted towards left with increase in λ .

$$\text{cost function} = 0 + (\text{slope})^2 \cdot \lambda$$

• Join me on LinkedIn for the latest updates on ML:

<https://www.linkedin.com/groups/7436898/>

- change θ value to create another best fit line.

$$\lambda \propto \frac{1}{\text{slope}}$$

→ inversely proportional.

$\lambda = 14$ make sure that our line doesn't overfit.

$\lambda \geq 0$, is a complexity parameter that controls the amount of shrinkage:

the larger the value of λ , the greater the amount of shrinkage.

The coefficient are shrink towards zero.

* θ value never becomes zero!

$$\begin{aligned} h_0(x) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3 \\ &= \theta_0 + 0.95x_1 + 0.82x_2 + \underline{\underline{0.10x_3}} \end{aligned}$$

It will get deleted.

\therefore Ridge Regression is used to introduce bias to the data in order to generalize the data and increase bias.

This is useful if you don't have much training data.

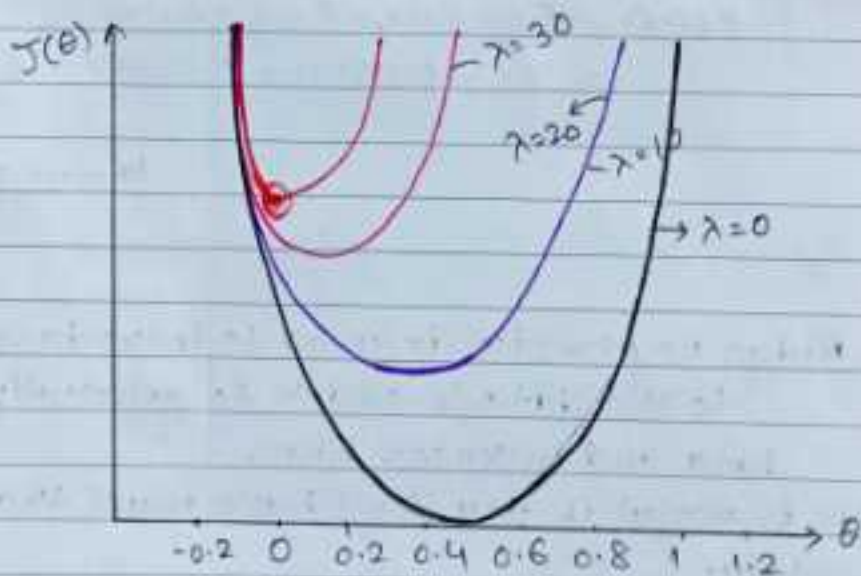
Lasso Regression

(L1 Regularization / L1 Norm)

- It is used to reduce the features. It helps in feature selection.

cost Function:

$$\text{cost function} = \frac{1}{n} \sum_{i=1}^n \left(h(\theta^{(i)}) - y^{(i)} \right)^2 + \lambda \sum_{i=1}^n |\text{slope}|$$



$$\begin{aligned} h(\theta) &= \theta_0 + \theta_1 \alpha_1 + \theta_2 \alpha_2 + \theta_3 \alpha_3 \\ &= \theta_0 + 0.54 \alpha_1 + 0.23 \alpha_2 + 0.10 \alpha_3 \end{aligned}$$

least correlated.

- If data has outliers \rightarrow use Ridge Regression.

Lasso = Least Absolute Shrinkage and Selection Operator Regression.

- Lasso Regression tends to eliminate the weights of the least important features by setting their weights to zero.

Elastic Net

\rightarrow combination of L1 and L2 Regularization.

$$\text{cost function} = \frac{1}{m} \sum_{i=1}^m \left\{ h(\mathbf{x})^i - y^{(i)} \right\}^2 + \underbrace{\lambda (\text{slope})^2}_{L2} + \underbrace{\lambda |\text{slope}|}_{L1}$$

can be changed to MAE, RMSE, MSE.

•• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

Notes taken from John Starmmer of Stat Quest

Youtube Videos

Regularization: Ridge (L2) Regression

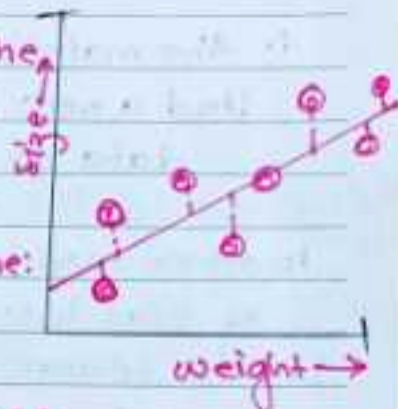
we find the line that results in the minimum sum of squared residuals.

∴ we end up with the eqn of the line:

$$\underline{\text{size}} = 0.95 + 0.75 \times \text{weight}$$

y-intercept

slope



when we have a lot of measurements, we can be fairly confident that least squares line accurately reflects the relationship between size and weight.



But what if we only have two measurements?

we fit new line. since the new line overlaps the two data points, the minimum sum of squared residuals = 0.

∴ New line eqn: $\text{size} = 0.4 + 1.3 \times \text{weight}$

sum of the squared residuals for testing data is large which means the new line has high variance.



in ML, new line (blue) is overfit to training data.

∴ the main idea behind Ridge Regression is to find a new line that doesn't fit the training data as well.

In other words, we introduce a small amount of Bias into how the new line is fit to the data but in return for that small amount of bias, we get a significant drop in variance.

∴ Ridge Regression can provide better long term predictions.

when least squares determines values for the parameters in this equation



$$\text{size} = \text{y-axis intercept} + \text{slope} \times \text{weight}$$

it minimizes....

the sum of the squared residuals.

• Download Machine Learning:
<https://t.me/AIMLDeepThought>

In contrast,

when Ridge Regression determines values for the parameters in this equation....



$$\text{size} = \text{y-axis intercept} + \text{slope} \times \text{weight}$$

... it minimizes
the sum of the squared residuals

+

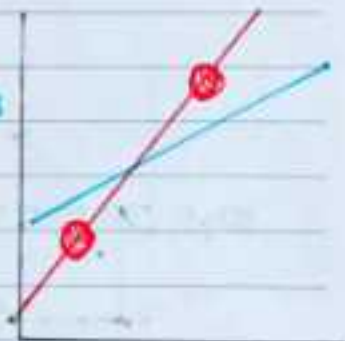
$\lambda \times \text{the slope}^2$

→ lambda.

this part adds a penalty to
the traditional least square
method.

and lambda(λ) determines how severe that
penalty is.

→ the sum of squared residuals
for the ~~residuals~~ least square fit
is 0 (because the line overlaps
the data points). and the slope
is 1.3.



$$\therefore 0 + \lambda \times (1.3)^2 = 0 + 1 \times (1.3)^2 = \underline{1.69}$$

for blue line →

slope = 0.8

$$(0.3)^2 + (0.1)^2$$

$$\lambda + (0.8)^2$$

λ = 1 all together

$$\underline{0.74}$$

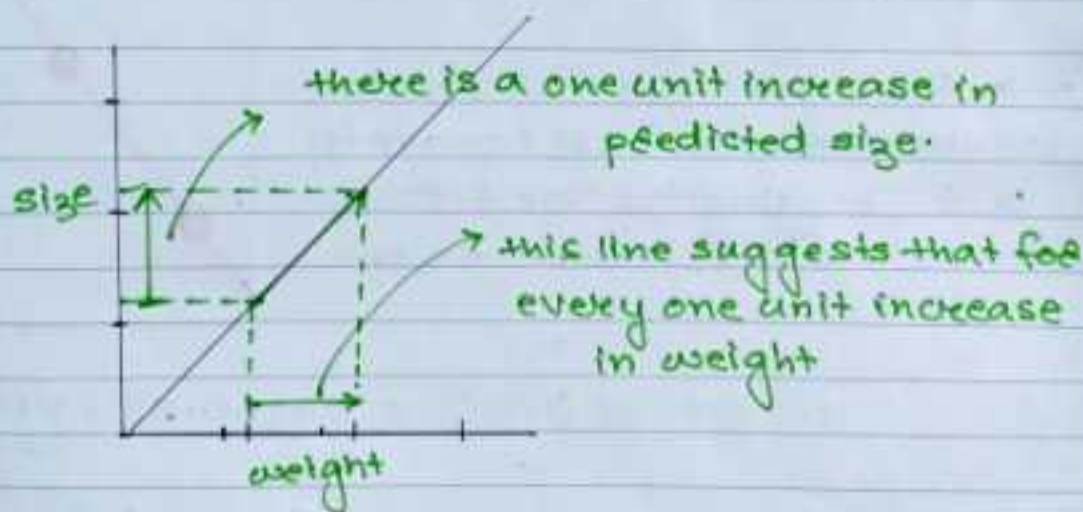
∴ Ridge Regression line :-

Red = 1.69, Blue = 0.74

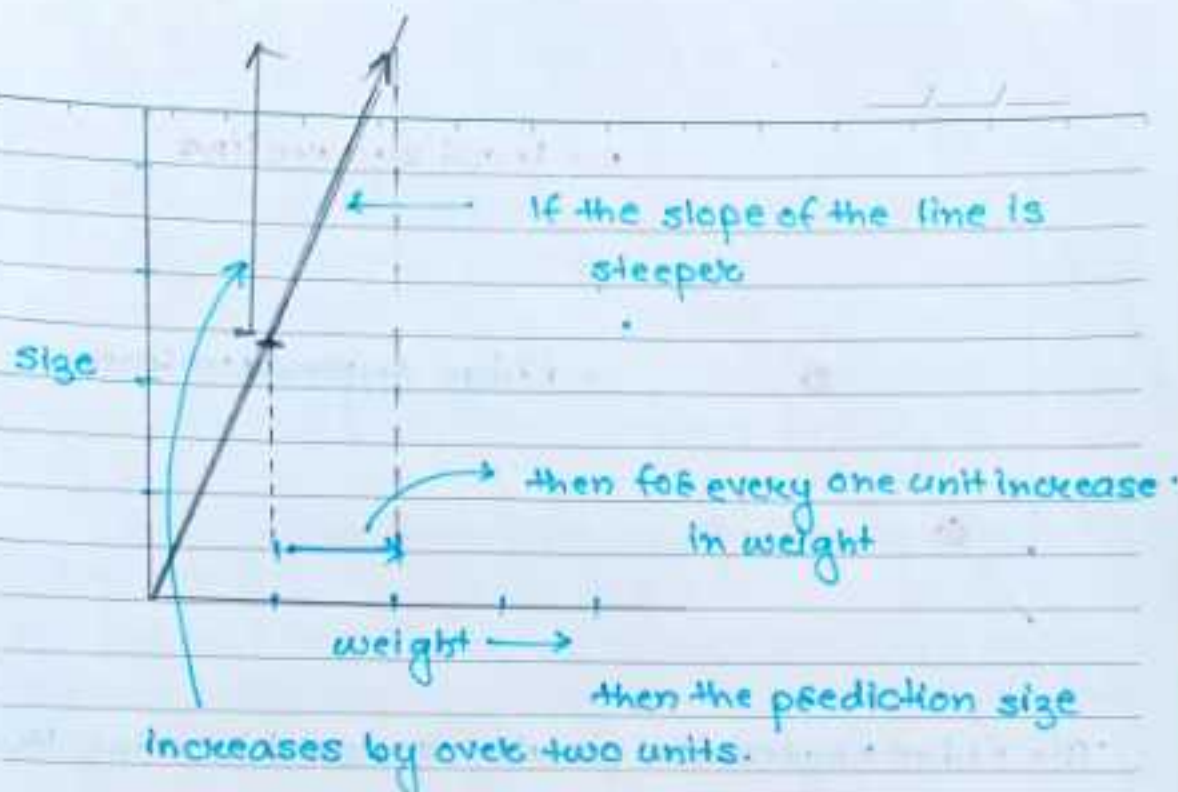
Thus, if we wanted to minimize the sum of the squared residuals plus the Ridge Regression penalty, we would choose the Ridge Regression Line over the least square line.

without the small amount of Bias that the penalty creates, the least squares fit has a large amount of Variance.

In contrast, the Ridge Regression line, which has the small amount of Bias due to the penalty, has less variance.



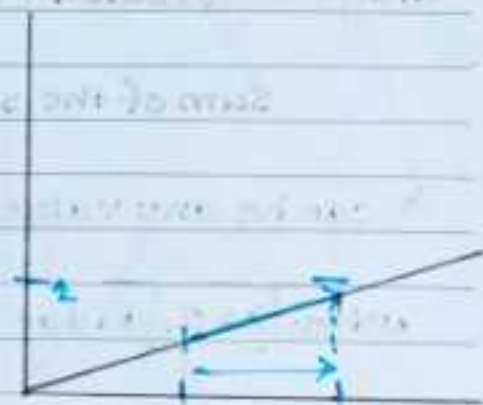
•• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>



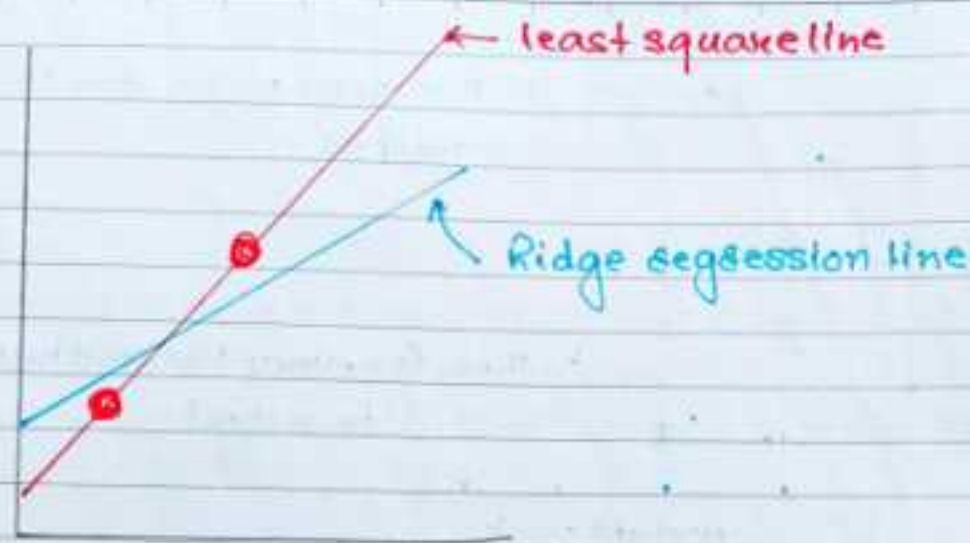
In other words, when the slope of the line is steep, then the prediction for size is very sensitive to relatively small changes in weight.

when the slope is small, then for every one unit increase in weight

the prediction for size barely increases.



In other words, when the slope of the line is small, then predictions for size are much less sensitive to changes in weight.



The Ridge Regression penalty resulted in a line that has a smaller slope ...

which means that predictions made with the Ridge Regression line are less sensitive to weight than the least square line.

Ridge Regression (RR)

Sum of the squared residuals + $\lambda \times (\text{slope})^2$

λ can be any value from 0 to positive infinity.

when $\lambda = 0$, Ridge Regression line = least square line

$\lambda = 1$ RR ended up with a smaller slope than the least square line.

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

and the larger we make λ , the slope gets asymptotically close to 0.



So, the larger λ gets, our predictions for size becomes less and less sensitive to weight.

so how do we decide what value to give λ ?

we just try a bunch of values for λ and use cross-validation, typically 10-fold cross validation, to determine which one results in the lowest variance.

... uptill now RR was for continuous variable.

However, RR also works when we use discrete variable.

Discrete variable:

y -intercept \rightarrow corresponds to the average size of the mice on the Normal diet.

$$\text{Size} = 1.5 + 0.7 \times \text{High fat disease}$$

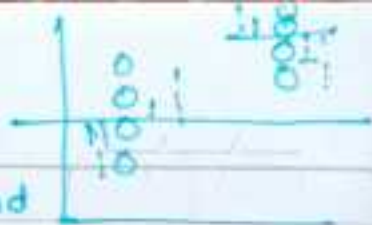
sum of these two is the prediction for the size of the mice on the High fat diet.

$$\text{size} = 1.5 + 0.7 \times \text{High fat diet}$$

y -intercept



these distance between the data and the means are minimized.



when RR determines value for the parameters in the equation...

... it minimizes \rightarrow

the sum of the squared residuals
+
 $\lambda \times (\text{diet difference})^2$

$\lambda = 0$, least squared-error = RR line

$\lambda = 1$, then only way to minimize the whole eqn is to shrink diet distance down.

In other words, as λ gets larger, our prediction for the size of the mice on the high fat diet becomes less sensitive to the difference between the Normal diet and High-fat diet.

The whole point of doing RR is because small sample size like these can lead to poor least squares estimates that result in terrible machine learning predictions.

•• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

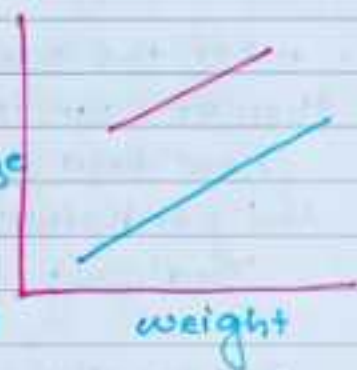
Ridge Regression can also be applied to Logistic Regression.

$$= \text{the sum of the likelihoods} + \lambda(\text{slope})^2$$

Note:- when applied to Logistic Regression, Ridge Regression optimizes the sum of the likelihoods instead of the squared residuals because Logistic Regression is solved using maximum likelihood.

" Ridge Regression helps reduce variance by shrinking parameters and making our predictions less sensitive to them.

In general, RR penalty contains all of the parameters size except for the y-intercept.



the sum of the squared residuals

$$+ \lambda(\text{slope}^2 + \text{dist distance}^2)$$



least squares can't find a single optimal solution, since any line that goes through the dot will minimize the sum of the squared residuals.

but RR can find a solution with cross validation and the RR penalty that favours smaller parameter values.



$$\begin{aligned} &\text{sum of the squared residuals} \\ &+ \\ &\lambda * (\text{slope})^2 \end{aligned}$$

Summary:

- when the sample sizes are relatively small, then RR can improve predictions made from new data (i.e. reduce variance) by making the predictions less sensitive to the Training data.

RR penalty itself is λ times the sum of all squared parameters, except for the y-intercept. and λ is determined using cross validation.

Lasso Regression: (L1)

$$\text{Ridge Regression Penalty} = \lambda \times (\text{slope})^2$$

* Lasso Regression \rightarrow

$$\begin{aligned} &\text{sum of all the squared residuals} \\ &+ \lambda \times |\text{slope}| \end{aligned}$$

Lasso Regression Penalty contains all of the estimated parameters except for the y-intercept.

\therefore Ridge and Lasso Regression shrink parameters, they don't have to shrink them all equally.

Big difference between Ridge and Lasso Regression is that Ridge Regression can only shrink the slope asymptotically close to 0 while Lasso Regression can shrink the slope all the way to 0.

LR can exclude useless variables from equations, better than RR at reducing the variance in models that contain a lot of useless variables.

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

Elastic Net Regression:

Elastic-Net Regression starts with least-squares then combines the Lasso Regression penalty with the Ridge Regression penalty.

$$\begin{aligned} & \text{sum of the squared residuals} \\ & + \\ & \lambda_1 \times | \text{variable}_1 | + \dots + | \text{variable}_n | \\ & + \\ & \lambda_2 \times (\text{variable}_1)^2 + \dots + (\text{variable}_n)^2 \end{aligned}$$

Note: LR and RR penalty get their own λ s.

The hybrid Elastic Net Regression is especially good at dealing with situations when there are correlations between parameters.

This is because on it's own, Lasso Regression tends to pick just one of the correlated terms and eliminates the others whereas RR tends to shrink all of the parameters for the correlated variables together.

•• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

By combining LR and RR,

Elastic-Net Regression groups and shrinks the parameters associated with the correlated variables and leaves them in eqn or removes them all at once.

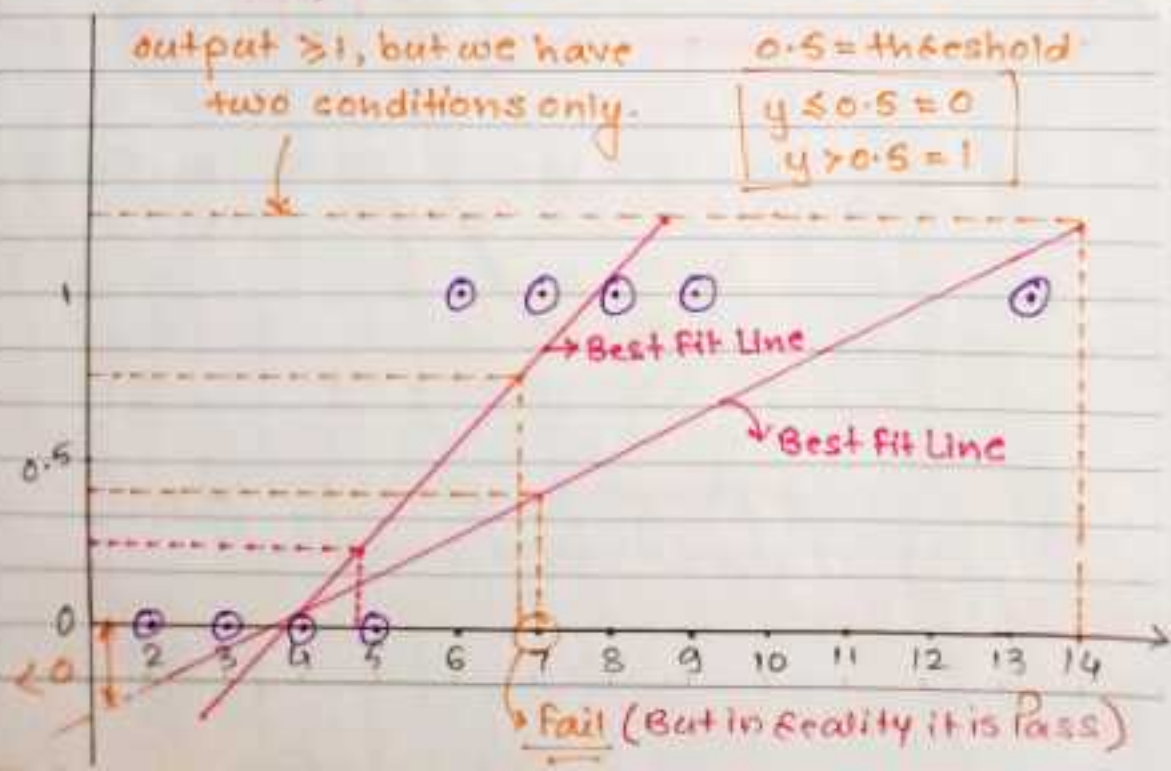
Logistic Regression

- classification problem

Eg. dataset:

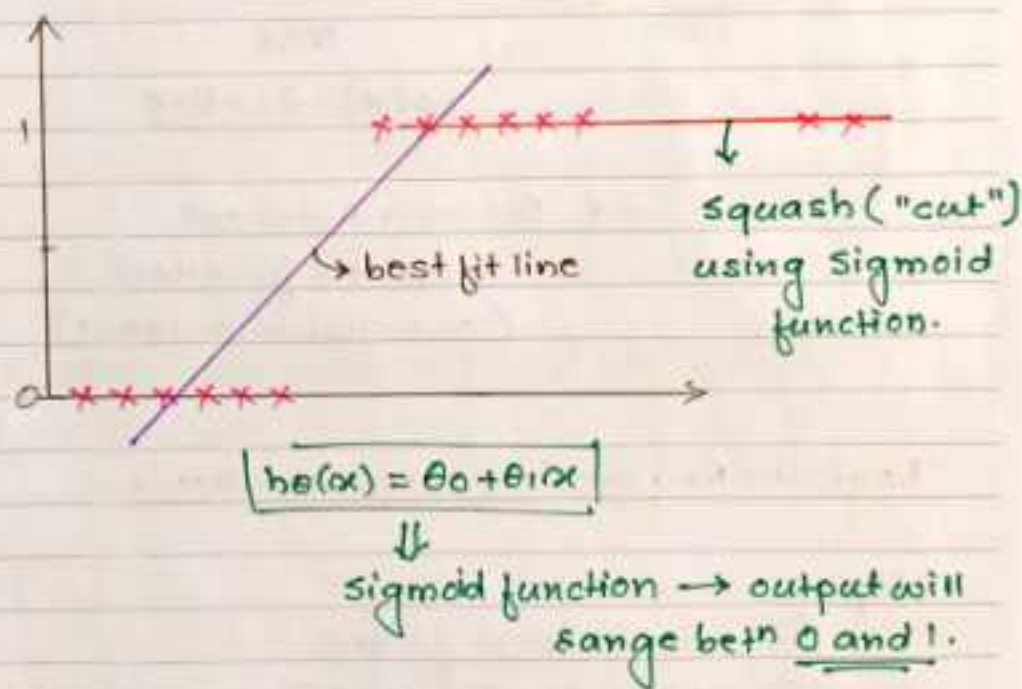
Study hours	Play hours	O/P (Pass/Fail)
2	8	Fail
3	7	Fail
6	3	Pass
outliers \rightarrow 1	4	Pass

* we cannot perform regression, we need to perform classification into Pass/Fail.



why we use logistic Regression when we can solve classification problem using Linear Regression?

- • due to outliers best fit line gets change and results will be wrong.
- we cannot remove outliers always.
- Threshold can't be changed, once fixed.
- In logistic we squash ("cut") the line. we will not change the line.



•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

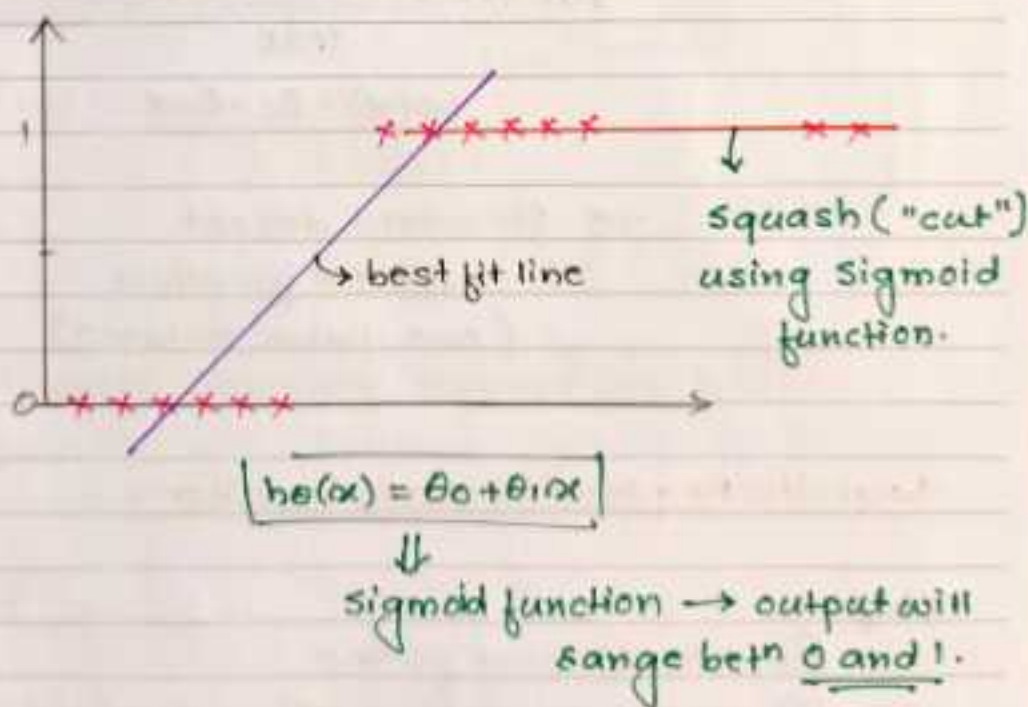
so, even if a person study for 6 Thrs
 $y < 0.5$

- ∴ the model will show Fail.
- ∴ we cannot use Linear Regression in this type of problem statement.

n we can
sing

change and

- we cannot remove outliers always.
- Threshold can't be changed, once fixed.
- In logistic we squash ("cut") the line. we will not change the line.



Sigmoid Function:

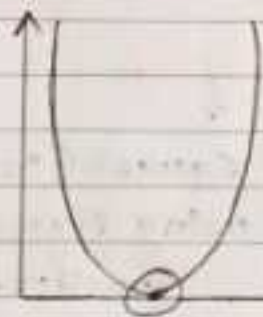
$$\text{sigmoid function} = \frac{1}{1 + e^{-x}} \rightarrow \text{bet'n } 0 \text{ \& } 1.$$

1. create a best fit line.
2. Squashing \rightarrow sigmoid function.

Linear Regression cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m \underbrace{\{h_{\theta}(\alpha)^{(i)} - y^{(i)}\}^2}_{\text{MSE}}$$

$$\therefore h_{\theta}(\alpha) = \theta_0 + \theta_1 \alpha$$



\rightarrow Gradient descent
convex function
(one global minima)

Logistic Regression cost function:

steps:

- ① create best fit line
- ② apply squashing using sigmoid function.

$$J(\theta_0, \theta_1) = \frac{1}{m} \sum_{i=1}^m \{h_{\theta}(x^{(i)}) - y^{(i)}\}^2$$

→ sigmoid function

$$h_{\theta}(x) = \sigma(\theta_0 + \theta_1 x)$$

let $z = \theta_0 + \theta_1 x$

let $z = \theta_0 + \theta_1 x$

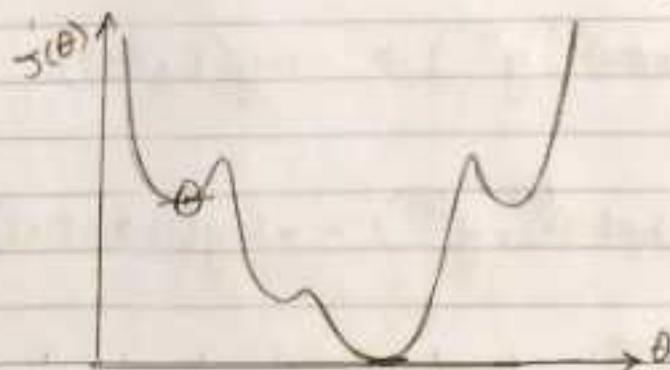
$$h_{\theta}(x) = \sigma(z)$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-z}}$$

$$\therefore h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

$$\therefore h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

but, after applying sigmoid function, cost function will become non-convex function and have a chances to get local minima.



← non-convex function.

- change the cost function to solve the convexity problem.

• Log Loss Cost Function

$$\text{cost function} = \begin{cases} -\log(h_{\theta}(x)), & \text{if } y=1 \\ -\log(1-h_{\theta}(x)), & \text{if } y=0 \end{cases}$$

→ convex function

$$h_{\theta}(x) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}$$

$$\text{cost}(h_{\theta}(x)^{(i)}, y^{(i)}) = -y \log(h_{\theta}(x)) - (1-y) \log(1-h_{\theta}(x))$$

→ This will never give local minima.

$$\text{if } y=1, \quad \text{cost}(h_{\theta}(x)^{(i)}, y^{(i)}) = -\log(h_{\theta}(x))$$

$$\text{if } y=0, \quad \text{cost}(h_{\theta}(x)^{(i)}, y^{(i)}) = -\log(1-h_{\theta}(x))$$

y = truth value

minimize cost function $J(\theta_0, \theta_1)$ by changing θ_0, θ_1 .

Convergence Algorithm:

Repeat Convergence

{

$$\theta_j := \theta_j - \alpha \cdot \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1)$$

}

for $j=0$ and $j=1$

- by default take threshold = 0.5
using ROC and TOC curve, we can
define threshold.

**** Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>**

**** Download Machine Learning:
<https://t.me/AIMLDeepThought>**

Performance Metrics

1. Confusion matrix
2. Accuracy
3. Precision
4. Recall
5. F-Beta Score



dataset

f_1	f_2	o/p	$P = \hat{y}$ (model prediction)
-	-	0	1
-	-	1	1
-	-	0	0
-	-	1	1
-	-	1	1
-	-	0	1
-	-	1	0

Confusion matrix

1 → correct prediction
0 → wrong prediction

		← Actual value (y)	
		1	0
↑ predicted value (y)	1	3	2
	0	1	1

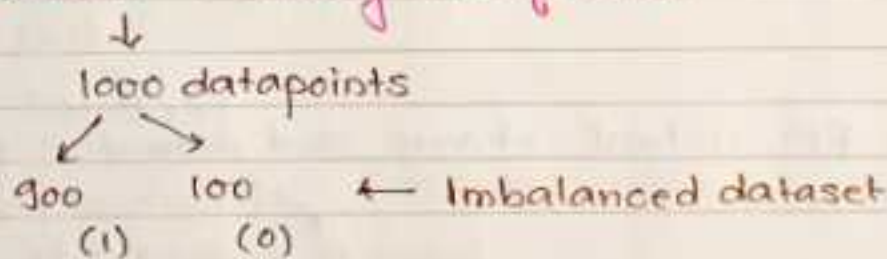
		← Actual value		
		1	0	(y)
↑ predicted value (\hat{y})	1	TP	FP	correct match } TP: True Positive TN: True Negative
	0	FN	TN	
				wrong match } FP: False Positive FN: False Negative

Accuracy

$$\text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}$$

Eg. $\text{accuracy} = \frac{3+1}{3+2+1+1} = \frac{4}{7} = 57\%$

dataset → Binary classification



dumb model → 1 \Rightarrow we get 90% accuracy.

- If the accuracy is 90%, it is not sufficient. model is not good. To overcome this problem we can use Recall and Precision.

• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

Precision:

$$\text{Precision} = \frac{TP}{TP + FP}$$

← Actual

	1	0
↑ Predicted	TP	FP
0	FN	TN

out of all the actual values
how many are correctly
predicted.

• our aim is to reduce False Positive (FP)

Eg. ① mail → spam or ham

TP: actual → spam and predicted → spam
(good)

FN: actual → spam and predicted → ham
(not good)

FP: actual → ham and predicted → spam
(mail is not spam) (critical problem)
focus on reducing FP.

② Diabetes or Not diabetes

Actual → diabetes and Predicted → not
diabetes

critical problem → reduce FN.

Recall:

$$\text{Recall} = \frac{TP}{TP + FN}$$

→ out of all the predicted values how many are currently predicted.

Eg. Tomorrow the stock market is going to crash.

Two points → consumers → FN ↓
→ companies → FP ↓
(can take certain decision)

1 ← Actual

	1	0
1	TP	FP
0	FN	TN

↑ Predicted

consumer

for FP companies can take action: company (action: sale shares and discounted price)

F-Beta Score

$$\text{F-beta score} = \frac{(1 + \beta^2) \text{Precision} \times \text{Recall}}{\beta^2 \cdot \text{Precision} + \text{Recall}}$$

- ① If FP and FN are both important
 $\beta = 1$

$$\therefore \text{F1 score} = \frac{2 \cdot P \times R}{P + R}$$

- ② If FP is more important than FN
 $\beta = 0.5$

$$\text{F-0.5 score} = \frac{(1 + 0.25) P \times R}{0.25 \times P + R}$$

- ③ If $\text{FN} \gg \text{FP}$, (FP is less important than FN),
 $\beta = 2$

$$\text{F2 score} = \frac{(1 + 4) P \times R}{4 \times P + R}$$

• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

• Download Machine Learning:
<https://t.me/AIMLDeepThought>

In-depth Insights

Logistic Regression

Classification

Classification problem, is just like the regression model problem, except that the values we now want to predict take on only a small number of discrete values.

$y \in \{0, 1\}$ 0: " -ve class"
Binary class 1: " +ve class"

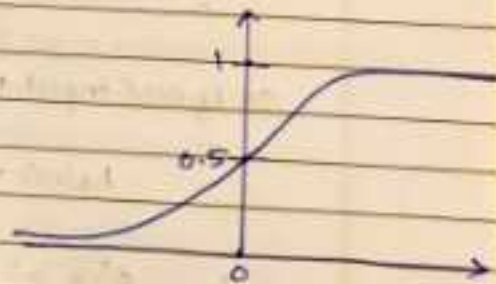
Hypothesis Representation

Intuitively, it doesn't make sense for $h(\alpha)$ to take values larger than 1 or smaller than 0 when we know that $y \in \{0, 1\}$. To fix this, let's change the form for our hypothesis $h(\alpha)$ to satisfy $0 \leq h(\alpha) \leq 1$.

This is accomplished by plugging $\theta^T \alpha$ into the logistic function.

Our new form uses the "sigmoid function" also called the "logistic function".

$$\begin{aligned} h(\alpha) &= g(\theta^T \alpha) \\ z &= \theta^T \alpha \\ g(z) &= \frac{1}{1 + e^{-z}} \end{aligned}$$



Function $g(z)$, shown here, maps any real number to the $(0, 1)$ interval, making it useful for transforming an arbitrary-valued function into a function better suited for classification.

$h(\alpha)$ will give us the probability that our output is 1.

Eg. if $h(\alpha) = 0.7 \rightarrow$ output: 1

thus $h(\alpha) = 0.3 \rightarrow$ 0

$$h_{\theta}(x) = P(y=1|x;\theta) = 1 - P(y=0|x;\theta)$$

$$P(y=0|x;\theta) + P(y=1|x;\theta) = 1$$

Decision Boundary

In order to get our discrete 0 or 1 classification, we can translate the output of the hypothesis function as follows:

$$h_{\theta}(x) \geq 0.5 \rightarrow y=1$$

$$h_{\theta}(x) < 0.5 \rightarrow y=0$$

The way our logistic function g behaves is that when its input is greater than or equal to zero, its output is greater than or equal to 0.5:

$$g(z) \geq 0.5 \text{ when } z \geq 0$$

Remember:

$$z=0, e^0=1 \Rightarrow g(x)=1/2$$

$$z \rightarrow \infty, e^{\infty} \rightarrow \infty \Rightarrow g(x)=1$$

$$z \rightarrow -\infty, e^{\infty} \rightarrow \infty \Rightarrow g(x)=0$$

so if our input to g is $\theta^T x$, then that means

$$h_{\theta}(x) = g(\theta^T x) \geq 0.5 \text{ when } \theta^T x \geq 0$$

$$\therefore \theta^T x \geq 0 \rightarrow y=1$$

$$\theta^T x < 0 \rightarrow y=0$$

The decision boundary is the line that separates the area where $y=0$ & $y=1$. It is created by our hypothesis function.

Eg.

$$\theta = \begin{bmatrix} 5 \\ -1 \\ 0 \end{bmatrix}$$

$$y = 1 \text{ if } 5 + (-1)x_1 + 0 \cdot x_2 \geq 0$$

$$5 - x_1 \geq 0$$

$$-x_1 \geq -5$$

$$\underline{x_1 \leq 5}$$

In this case, our decision boundary is a straight vertical line placed on the graph where $x_1 = 5$, and everything to the left of that denotes $y = 1$, while everything to the right denotes $y = 0$.

again, the input to the sigmoid function $g(z)$ (e.g. $\theta^T x$) doesn't need to be linear, could be a function that describes a circle (e.g. $z = \theta_0 + \theta_1 x_1^2 + \theta_2 x_2^2$) or any shape to fit our data.

Cost Function

we cannot use the same cost function that we use for linear regression because the Logistic Function will cause the output to be wavy, causing many local optima:

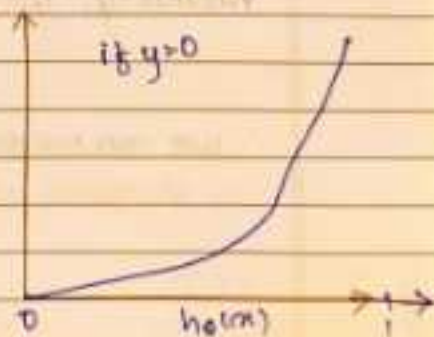
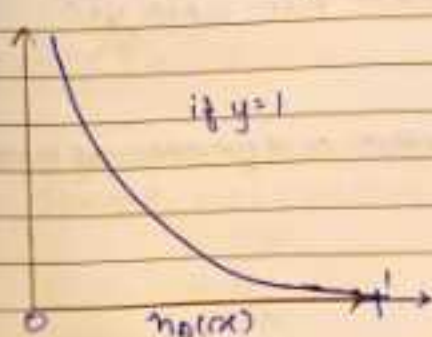
In other words, it will not be a convex function.

Cost function for logistic function:

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \text{cost}(h(\theta^T x^{(i)}), y^{(i)})$$

$$\text{cost}(h(\theta^T x), y) = -\log(h(\theta^T x)) \text{ if } y = 1$$

$$\text{cost}(h(\theta^T x), y) = -\log(1 - h(\theta^T x)) \text{ if } y = 0$$



$$\text{cost}(h_0(x), y) = 0 \text{ if } h_0(x) = y$$

$$\text{cost}(h_0(x), y) \rightarrow \infty \text{ if } y=0 \text{ and } h_0(x) \rightarrow 1$$

$$\text{cost}(h_0(x), y) \rightarrow \infty \text{ if } y=1 \text{ and } h_0(x) \rightarrow 0$$

If our correct answer 'y' is 0, then the cost function will be 0 if our hypothesis function also outputs 0. If our hypothesis approaches 1, then the cost function will approach infinity.

If our correct answer 'y' is 1, then the cost function will be 0 if our hypothesis function outputs 1. If our hypothesis approaches 0, then the cost function will approach infinity.

Note that writing the cost function in this way guarantees that $J(\theta)$ is convex for logistic function regression.

Simplified Cost function and Gradient Descent

We can compress our cost function's two conditional cases into one case.

$$\text{cost}(h_0(x), y) = -y \cdot \log(h_0(x)) - (1-y) \cdot \log(1-h_0(x))$$

Thus, when we substitute $y=1$ in abv eqn,

$$\text{cost}(h_0(x), y) = -y \cdot \log(h_0(x))$$

similarly, we get another term when $y=0$.

\therefore we can write our entire cost function as follows:

$$J(\theta) = -\frac{1}{m} \sum_{i=1}^m \left[y^{(i)} \cdot \log(h(\theta^T x^{(i)})) + (1 - y^{(i)}) \cdot \log(1 - h(\theta^T x^{(i)})) \right]$$

A vectorized implementation is:

$$h = g(X\theta)$$

$$J(\theta) = -\frac{1}{m} \left(-y^T \log(h) - (1-y)^T \log(1-h) \right)$$

Gradient Descent:

General form of gradient descent is:

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta) \end{array} \right.$$

we can work out the derivative part using calculus to get:

$$\text{Repeat } \left\{ \begin{array}{l} \theta_j := \theta_j - \frac{\alpha}{m} \sum_{i=1}^m (h(\theta^T x^{(i)}) - y^{(i)}) \cdot x_j^{(i)} \end{array} \right.$$

Notice that this algorithm is identical to the one we used in linear regression. we still have to simultaneously update all values in theta. ($h(\theta^T x)$ is different)

A vectorized implementation is:

$$\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \bar{y})$$

Advanced Optimization

cost function $J(\theta)$. want min $J(\theta)$.

Given θ , we have code that can compute

$$-J(\theta)$$

$$-\frac{\partial}{\partial \theta_j} J(\theta) \quad \text{for } j = 0, 1, \dots, m$$

Optimization algorithms:

- Gradient descent
- Conjugate gradient
- BFGS
- L-BFGS

advantages:- (1) No need to manually pick α .
(2) often faster than gradient descent

disadv:- more complex to implement.

Multiclass Classifications

Now we will approach the classification of data when we have more than two categories.

Instead of $y \in \{0, 1\}$ we will expand our definition so that $y \in \{0, 1, \dots, n\}$

since $y \in \{0, 1, \dots, n\}$, we divide our problem into $n+1$ ($+1$ because the index starts at 0) binary classification problems; in each one, we predict the probability that 'y' is a member of one of our classes.

$$y \in \{0, 1, \dots, n\}$$

$$h_0^0(\alpha) = P(y=0|\alpha; \theta)$$

$$h_0^1(\alpha) = P(y=1|\alpha; \theta)$$

⋮

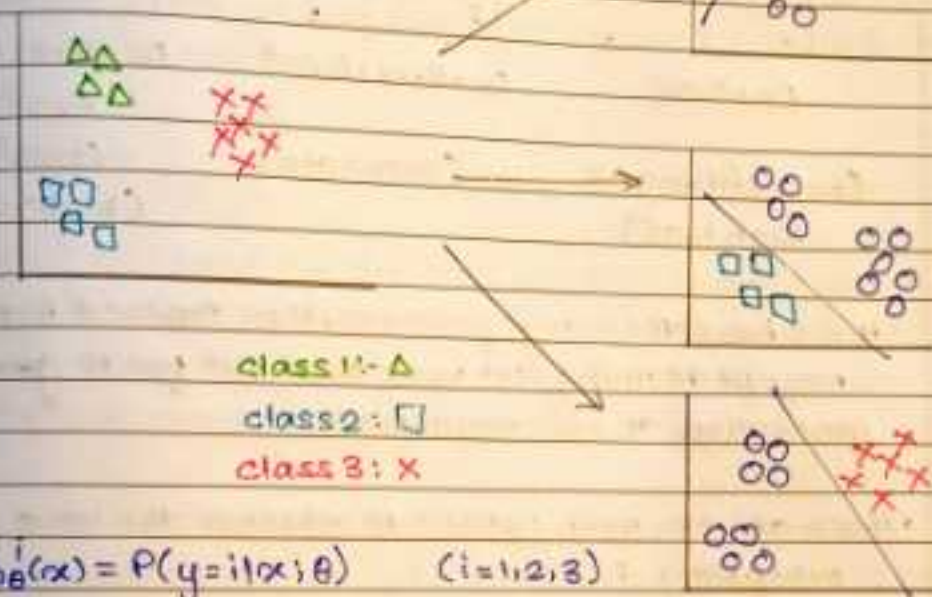
$$h_0^n(\alpha) = P(y=n|\alpha; \theta)$$

$$\text{prediction} = \max_k (h_0^k(\alpha))$$

we are basically choosing one class and then lumping all the others into a single second class. we do this repeatedly;

applying binary logistic regression to each case, and thus use the hypothesis that returned the highest value as our prediction.

one-vs-all (one-vs-rest)



To summarise:

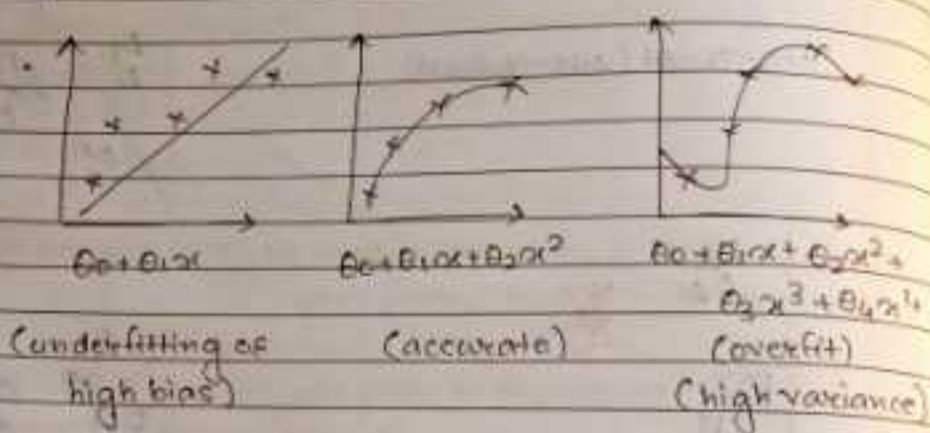
Train a logistic regression classifier, $h_{\theta}^i(x)$ for each class i to predict the probability that $y=i$

on a new input x , to make a prediction, pick the class that maximizes

$$\max_i h_{\theta}^{(i)}(x)$$

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

Regularization



If we have too many features, then learned hypothesis may fit training set very well but fail to generalize to new examples.

These are two main options to address the issue of overfitting:

1. Reduce the number of features:
 - manually select which features to keep.
 - use a model selection algorithm
2. Regularization
 - keep all the features, but reduce the magnitude of parameters θ_j .
 - Regularization works well when we have a lot of slightly useful ~~noisy~~ features.

Cost function

If we have overfitting from our hypothesis function, we can reduce the weight that some of the terms in our

function slowly by increasing their cost.
e.g.

We wanted to make the following function more quadratic:

$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

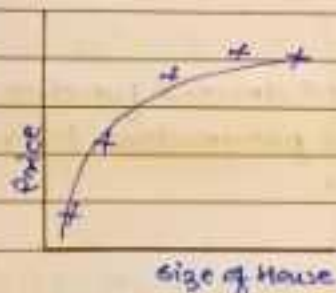
we will

without actually getting rid of these features or changing the form of our hypothesis, we can instead modify our cost function.

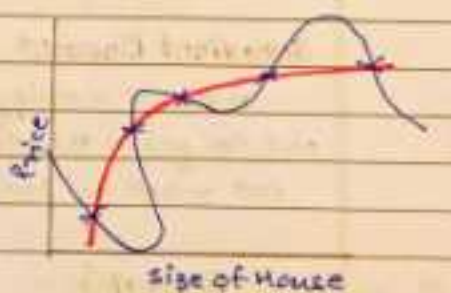
$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + 1000 \cdot \theta_3^2 + 1000 \cdot \theta_4^2$$

Now, in order for the cost function to get close to zero, we will have to reduce the values of θ_3 and θ_4 to near zero. This will in turn greatly reduce the values of $\theta_3 x^3$ and $\theta_4 x^4$ in our hypothesis function.

As a result, we see that new hypothesis looks like a quadratic function but fits the data better, due to extra small terms $\theta_3 x^3$ and $\theta_4 x^4$.



$$\theta_0 + \theta_1 x + \theta_2 x^2$$



$$\theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$$

Suppose we penalize and make θ_3, θ_4 really small.

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 + 1000 \theta_3^2 + 1000 \theta_4^2$$

pink graph shows this complete equation.

we could also regularize all of our theta parameters in a single summation as:

$$\min_{\theta} \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 + \lambda \sum_{j=1}^n \theta_j^2$$

The λ , or lambda, is the **Regularization parameter**. It determines how much the costs of our theta parameters are inflated.

using the above cost function with the earlier summation, we can smooth the output of our hypothesis function to reduce overfitting.

If lambda is chosen to be too large, it may smooth out the function too much and cause underfitting.

Regularized Linear Regression

Gradient Descent

we will modify gradient descent function to separate out θ_0 from the rest of the parameters because we do not want to penalize θ_0 .

Repeat $\{$

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \cdot x_{0i}$$

$$\theta_j := \theta_j - \alpha \left[\left(\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \cdot x_{ji} \right) + \right.$$

$$\left. \frac{\lambda}{m} \theta_j \right]$$

$j \in \{1, 2, \dots, n\}$

The term $\frac{\lambda}{m} \theta_j$ performs our regularization. With some manipulation our update rule can also be represented as:

$$\theta_j := \theta_j \left(1 - \alpha \frac{\lambda}{m}\right) - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i) \cdot x_j$$

\hookrightarrow always will be less than 1

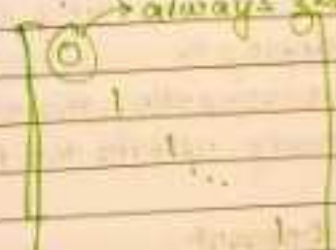
Intuitively we can see it as reducing the value of θ_j by some amount on every update.

2nd term is exactly same as it was before.

Normal Equation:

To add in regularization, the equation is the same as our original, except that we add another term inside the parentheses:

$$\theta = (X^T X + \lambda \cdot L)^{-1} X^T y$$

where $L =$ 

L is a matrix with 0 at the top left and 1's down the diagonal, with 0's everywhere else.

It should have dimension $(n+1) \times (n+1)$. Intuitively, this is the identity \rightarrow identity matrix (though we are not including 0's), multiplied with a single real no. λ .

If $m < n$, then $X^T X$ is non-invertible. However, when we add the term $\lambda \cdot L$, then $X^T X + \lambda \cdot L$ becomes invertible.

Regularized Logistic Regression

cost function:

cost function for logistic regression was:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(h_{\theta}(x^i)) + (1-y^i) \cdot \log(1-h_{\theta}(x^i))]$$

we can regularize this equation by adding a term to the end:

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m [y^i \cdot \log(h_{\theta}(x^i)) + (1-y^i) \cdot \log(1-h_{\theta}(x^i))] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2$$

means to explicitly exclude the bias term, θ_0 , i.e. the θ vector is indexed from 0 to n (holding $n+1$ values, θ_0 through θ_n).

and this sum explicitly skips θ_0 , by running from 1 to n , skipping 0.

Thus, when computing the equation, we should continuously update the two following equations:

Gradient Descent

repeat {

$$\theta_0 := \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) \cdot x_0^i$$

$$\theta_j := \theta_j - \alpha \left[\frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x_j^i + \frac{\lambda}{m} \theta_j \right]$$

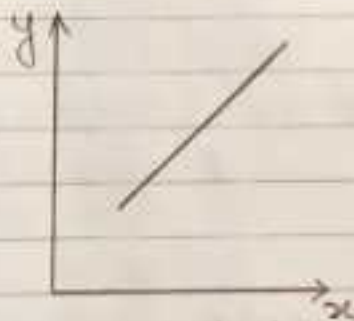
$\frac{\partial}{\partial \theta_j} J(\theta)$ $j = (1, 2, \dots, n)$

Support Vector Machine (SVM)

It can solve both classification and Regression problem.

1. classification \rightarrow SVC (Support Vector classifier)
2. Regression \rightarrow SVR (Support Vector Regressor)

some basics:



Equation of line:

$$y = mx + c \quad \text{OR}$$

$$y = \beta_0 + \beta_1 x \quad \text{OR}$$

$$ax + by + c = 0$$

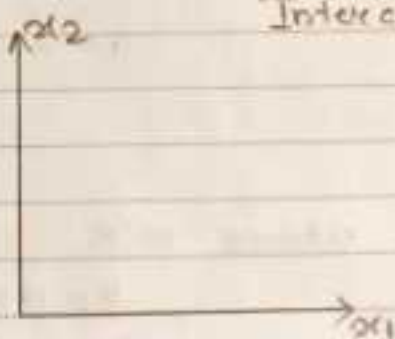
$$\therefore y = \left(-\frac{a}{b} \right) x - \left(\frac{c}{b} \right)$$

coefficient
Intercept

$$ax_1 + bx_2 + c = 0$$
$$w_1 x_1 + w_2 x_2 + b = 0$$

$$\therefore w^T x + b = 0$$

If line passes through origin $\therefore \underline{w^T x = 0}$

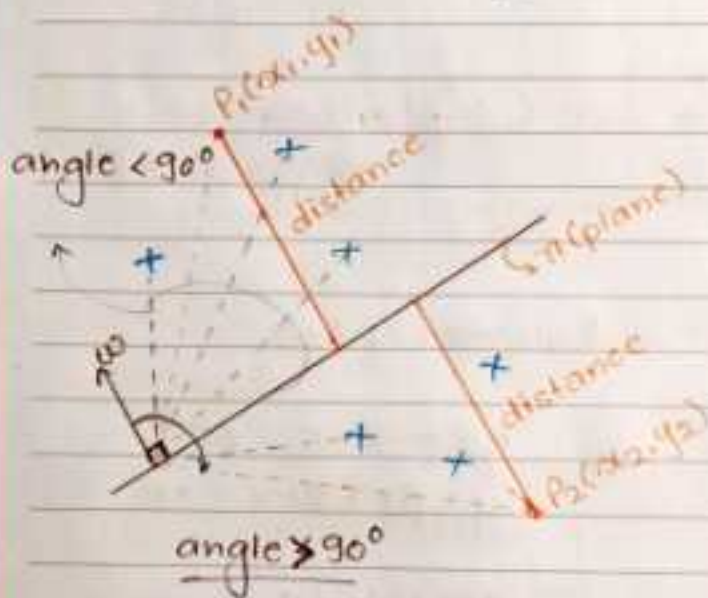


matrix multiplication:

$$\begin{bmatrix} w_1 \\ w_2 \end{bmatrix} \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$w^T x + b = 0. \quad (w^T x: w^T \text{ Transpose } x)$$

Equation of line passing through origin is:
 $w^T x = 0$



we have to find the distance of point from the plane.

(n = line in 2D and plane in 3D)

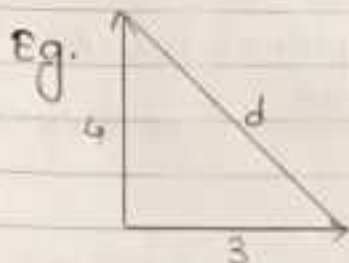
Distance of a point to the plane,

$$\text{distance}(d) = \frac{w^T P_1}{\|w\|}$$

$$= \|w\| \cdot \|P\| \cdot \cos \theta$$

where $w^T P_1$: w Transpose P_1 , w : vector,
 $\|w\|$: magnitude of w

unit vector: A vector which has a magnitude of 1 is basically called unit vector.



Now,

$$d = \sqrt{3^2 + 4^2}$$
$$= \sqrt{25}$$
$$\therefore d = 5$$

where vector, $\hat{d} = \frac{d}{\|d\|}$ → magnitude

$$\left(\frac{3}{5}, \frac{4}{5}\right) = \hat{d} = \frac{1}{5} \sqrt{(3)^2 + (4)^2} = \frac{1}{5} \sqrt{25}$$
$$= 1.$$

∴ unit vector is a way to get focused on direction not on magnitude.

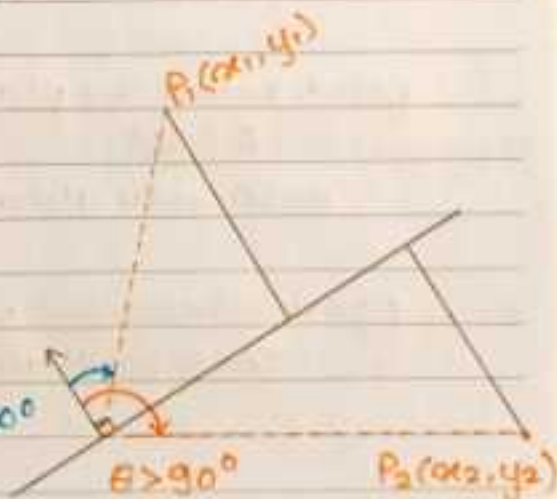
upward vector →

$$d = \frac{\omega^T P_i}{\|\omega\|}$$
$$d = \|\omega\| \cdot \|P_i\| \cdot \cos\theta$$

point above the plane
as $\theta < 90^\circ$

$\cos\theta$ will always be +ve.

$\theta < 90^\circ$

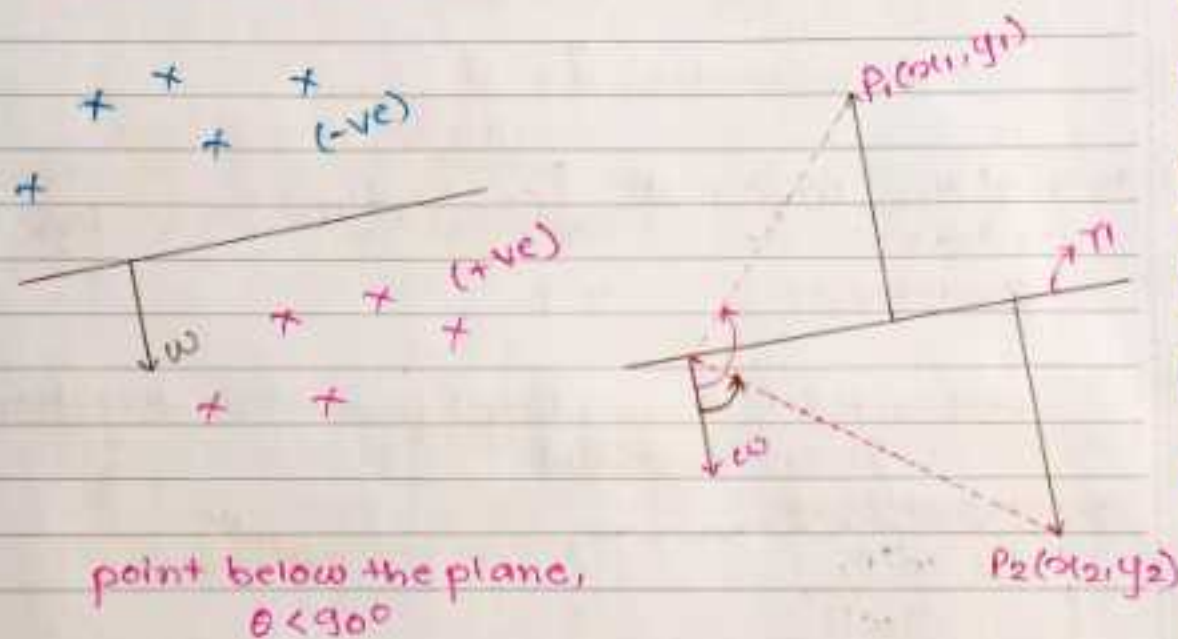


point below the plane, as $\theta > 90^\circ$
 $\cos\theta$ will always be -ve.

• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

- If any point falling above the plane, then θ must be less than 90° (+value)
- If any point falling below the plane, then θ must be greater than 90° (-ve value)

Downward vector:



point below the plane,
 $\theta < 90^\circ$

$\cos \theta$ will always be '+ve'.

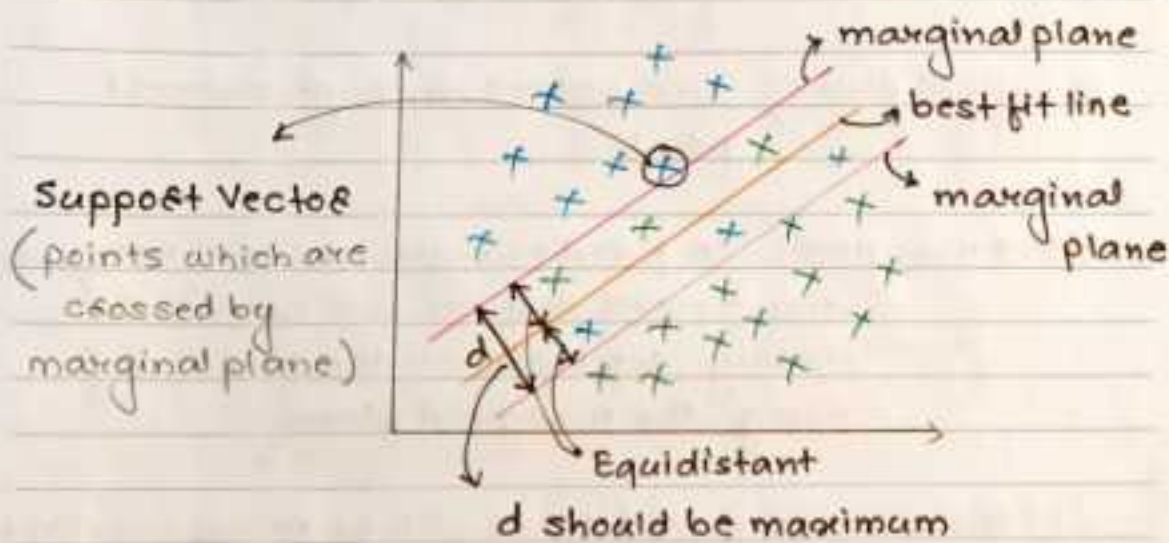
point above the plane, $\theta > 90^\circ$

$\cos \theta$ will always be '-ve'.

• Download Machine Learning:
<https://t.me/AIMLDeepThought>

Geometric Intuition Behind Support Vector Machine

Support Vector classifier (SVC)



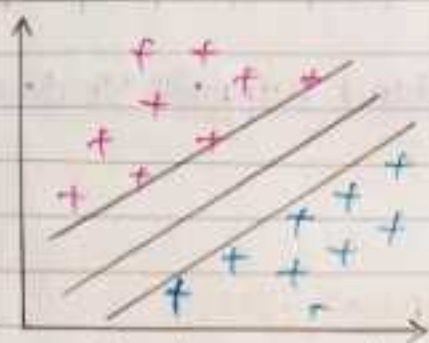
* you can have more than one support vector.

There will be many possible hyperplanes that separate different classes.

we have learnt in LR that the probability of a point belonging to any class given at very close to the hyperplane will be close to 0.5.

So, we want a hyperplane that separates (+ve) pts and (-ve) pts as far away as possible.

key idea of SVM: such hyperplane is called margin-maximizing plane.



Hard marginal



Soft marginal.

marginal plane: Line passes through nearest points.

Hard marginal: In hard margin, we will not find any errors and we will be able to clearly separate all the points by using the marginal plane.

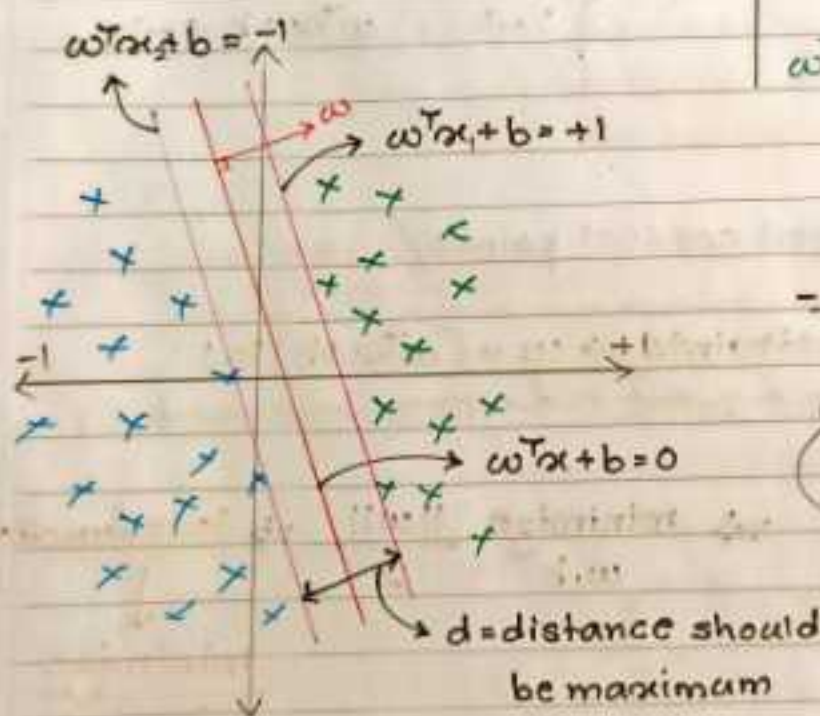
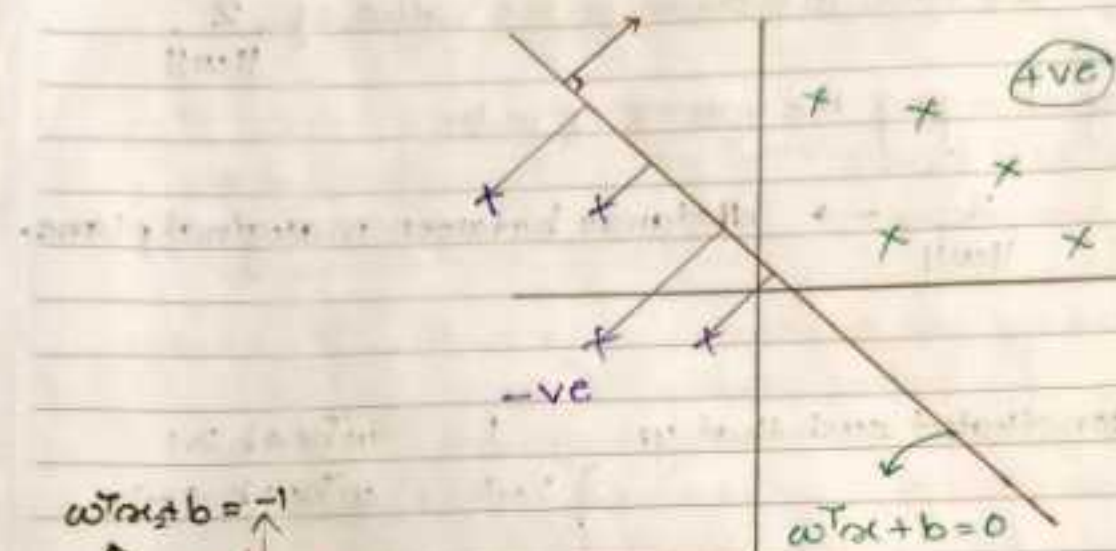
But, in real world there will be many overlapping with many errors, so marginal plane lines will be called soft margin.

→ marginal plane should be equidistance from best fit line.

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

• Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

SVM Mathematical Intuition:



$$\begin{aligned}\omega^T x_1 + b &= +1 \\ -\omega^T x_2 + b &= -1 \\ \omega^T (x_1 - x_2) &= 2\end{aligned}$$

∴ unit vector of ω , direction

$$\begin{aligned}\frac{\omega^T (x_1 - x_2)}{\|\omega\|} \\ &= \frac{2}{\|\omega\|}\end{aligned}$$

$$\therefore \text{distance } (d) = \frac{\omega^T (x_1 - x_2)}{\|\omega\|} = \frac{2}{\|\omega\|}$$

Cost Function:

we have to maximize the value of $\frac{2}{\|w\|}$

by changing the values of w, b .

$\frac{2}{\|w\|} \rightarrow$ distance between marginal plane.

constraint such that $y_i \begin{cases} 1 & w^T x + b \geq 1 \\ -1 & w^T x + b \leq -1 \end{cases}$

conditions

For all classified correct points,

constraints $\rightarrow y_i \times (w^T x + b) \geq 1$

maximize $\frac{2}{\|w\|} \Rightarrow$ minimize $\frac{\|w\|}{2} \Rightarrow$ loss function
 \downarrow
minimize

* loss function focus on minimization.

** Join me on LinkedIn for the latest updates on ML:
<https://www.linkedin.com/groups/7436898/>

Cost function:

minimize $\frac{\|w\|^2}{2}$ by changing w, b .

$$\min \frac{\|w\|^2}{2} + C_1 \sum_{i=1}^n \xi_i$$

Hinge loss for soft margin.

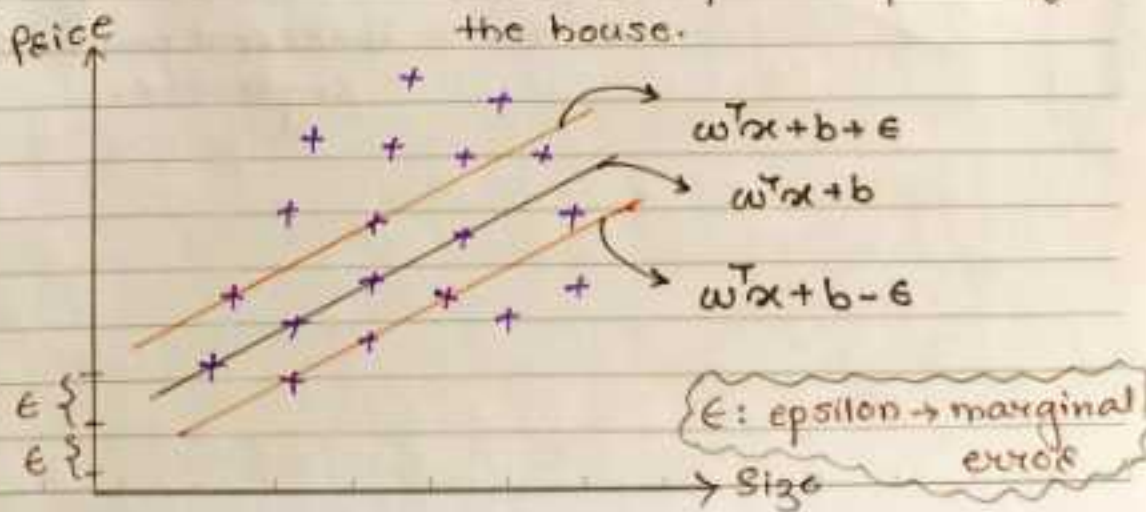
where, C_1 : How many points we can ignore for mis-classification.

Hyperparameter

ξ_i (Eta): Summation of the distance of incorrect data points from the marginal plane.

Support Vector Regressor:

Problem Statement: Based on the size of the house, we have to predict price of the house.



cost Function:

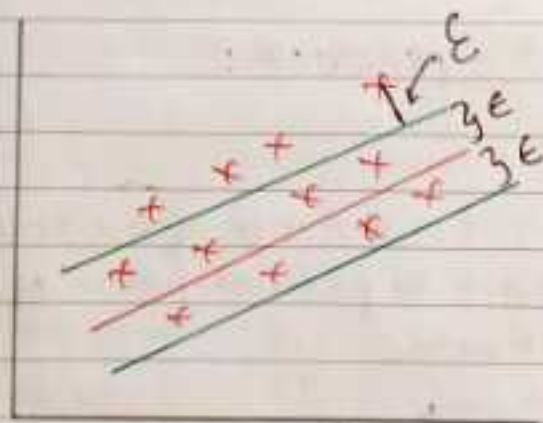
$$\underset{w, b}{\text{minimize}} \quad \frac{\|w\|^2}{2} + \left\{ C_i \sum_{i=1}^n \xi_i \right\} \rightarrow \text{Hinge loss}$$

MAE

$$\text{constraint: } |y_i - w_i x_i| \leq \epsilon + \xi_i$$

ϵ is margin
Truth point
predicted point
epsilon

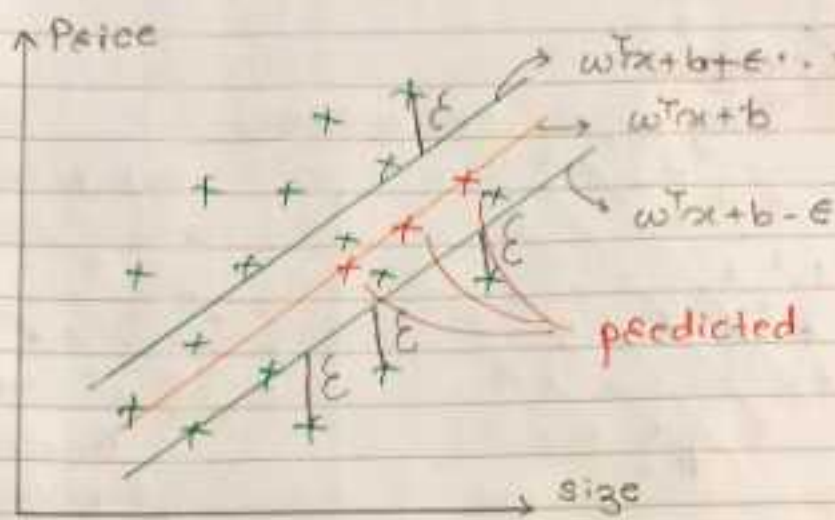
ϵ : margin of error (to decide original plane)
 ξ_i : error above the margin.



Hyperparameter:

- keep adjusting ϵ to get best margin
- we can't say incorrect point in regressor.

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>



In Regressor, no complete incorrect value, because it will be continuous value.

Q.

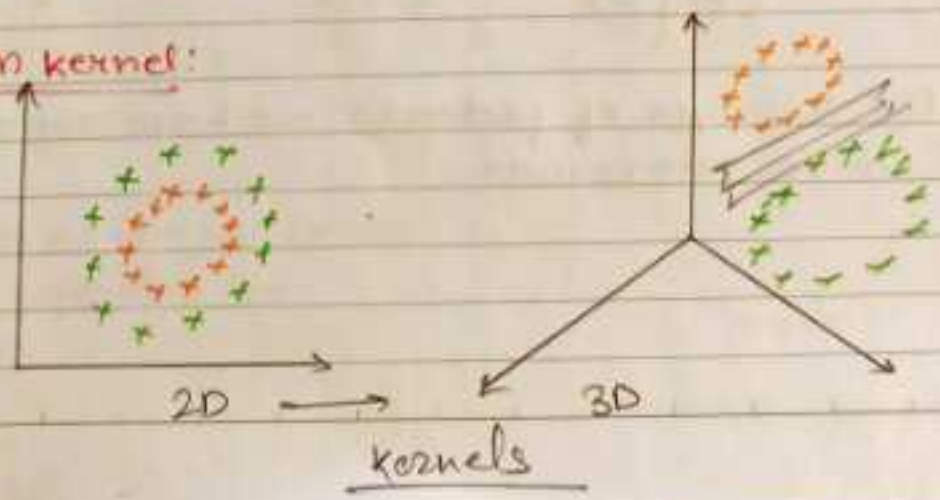
Is svm impacted by the outliers?

Yes, svm is impacted by the outliers.

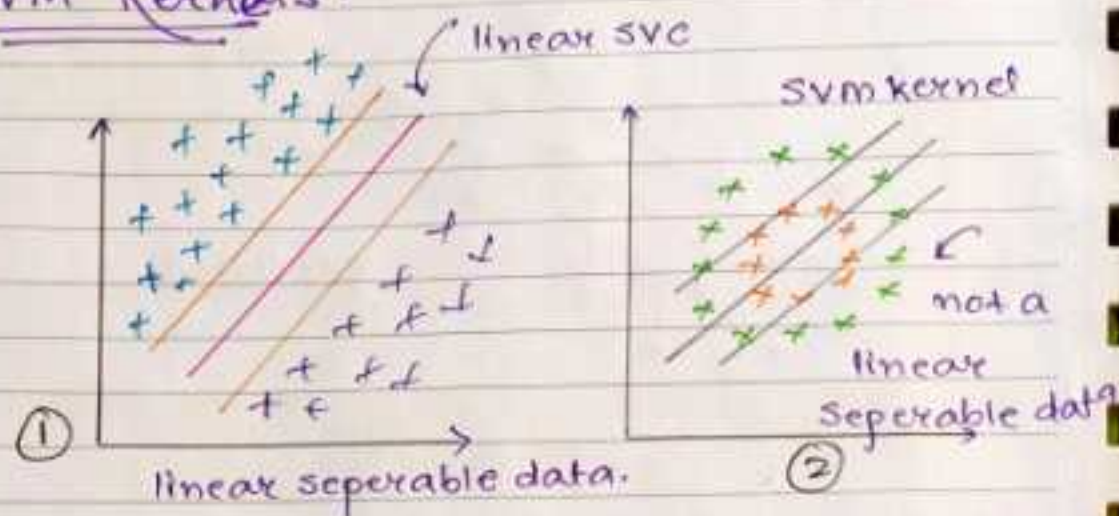
Does Standardization is need in svm?

Yes, we need to perform Normalization and standardization.

SVM kernel:



SVM Kernels:



- when we create this (1) type of best fit line and marginal plane, we are actually solving the linear separable data.
- → called as linear SVC. (Fig 1)
- If data is not a linear separable data, you will not be able to create best fit line and not able to create a marginal plane even though we create it, the accuracy will be very low. (Fig 2)
- For this type of problems, we have some more SVM kernels.

• Download Machine Learning:
<https://t.me/AIMLDeepThought>

what does SVM kernels do?

→ The main aim is to apply some transformation technique. (some mathematical formula) on the dataset.

This transformation increases the dimension of the data.

(mathematical formula)

SVM kernels → Transformation → Increasing the dimension of the data

linearly separable line

→ will give error

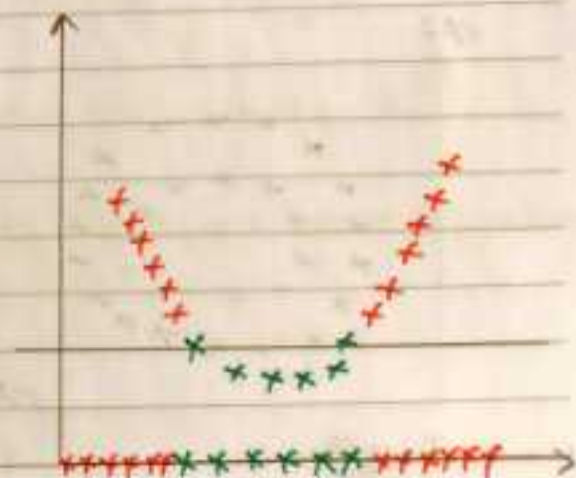


∴ we will transform the data from 1D → 2D

$$y = x^2$$

After →

Now, we can use linearly separable line.



$$y = x^2$$

so if $x = -7$ $y = 49$
 $x = -3$ $y = 9$ and so on.

what is the advantage of doing this transformation?

→ after transformation, we can apply linear SVM or SVC.

* when we divide convert 1D → 2D then we can divide all the points using single line which is called linear SVC.

Types of SVM kernel:

1. Polynomial kernel
2. RBF kernel
3. Sigmoid kernel

1. Polynomial kernel

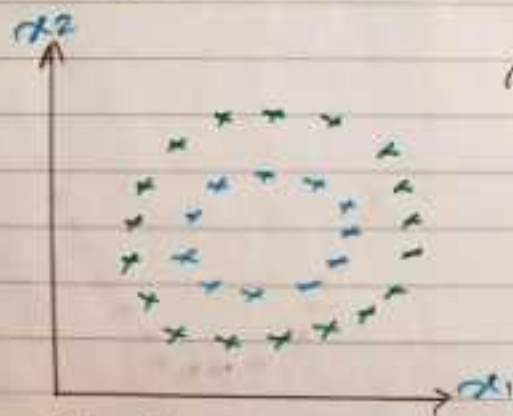
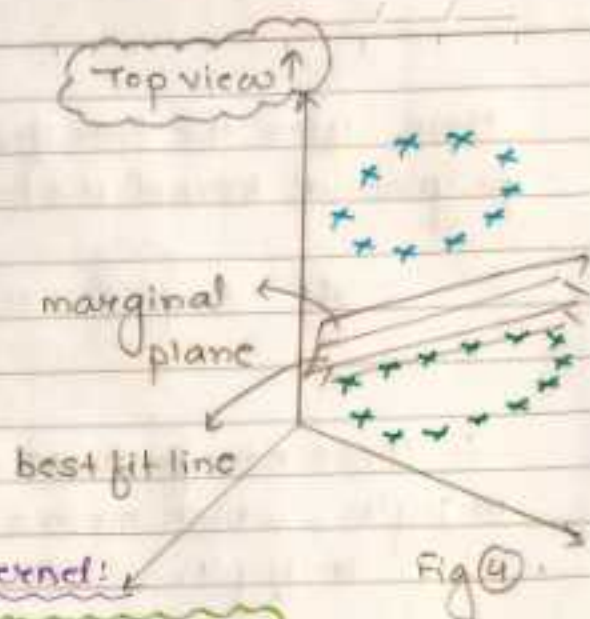


Fig 3.

→ not separable with best fit line.

so, we need to
convert 2D → 3D.

One main aim was to increase dimension,
 $2D \rightarrow 3D$
 so, hyperplane is created.



Formula for Polynomial kernel:

$$f(x_1, x_2) = (x_1^T x_2 + 1)^d$$

$d = \text{dimension}$

If we are converting $2D \rightarrow 3D$, the value of $d=3$.

$$\therefore x_1^T x_2 = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \cdot \begin{bmatrix} x_1 & x_2 \end{bmatrix}$$

$$= \begin{bmatrix} x_1^2 & x_1 \cdot x_2 \\ x_1 \cdot x_2 & x_2^2 \end{bmatrix}$$

3 unique values: $x_1^2, x_1 \cdot x_2, x_2^2$

Now, initially at the time of Fig 3, we have 3 features
 $x_1 \ x_2 \ y$

Join me on LinkedIn for the latest updates on Machine Learning: <https://www.linkedin.com/groups/7436898/>

Now, after the transformation / formula of polynomial kernel we have 6 features

$$x_1 \quad x_2 \quad \{x_1^2 \quad x_1 x_2 \quad x_2^2\} \quad y$$

∴ These new features can be plotted as the 3D in fig 4, which means that

$$\begin{aligned} \text{In fig 4, } x_1 &\text{ will be } x_1^2 \\ x_2 &\text{ will be } x_2^2 \\ z &\text{ will be } x_1 x_2 \end{aligned}$$

and once we have all these points, we will be able to clearly separate the points.

→ use polynomial kernel, to get better accuracy.

② Radial Basis Function Kernel (RBF kernel)

$$k(\vec{x}, \vec{x}_i) = e^{-\frac{\|\vec{x} - \vec{x}_i\|^2}{2\sigma^2}}$$

hyperparameter

③ Sigmoid Kernel

It can be used as the proxy for neural networks.

$$k(x, x_i) = \tanh(\sigma x^T x_i + c)$$

Mathematical formulation of SVM.

Objective: To find a hyper-plane that does margin maximization.

we need to find a hyperplane Π :- margin-maximization ...

Π :- $w^T x + b = 0$ (w not necessarily unit vector).

$$\Pi_+ : w^T x + b = 1$$

$$\Pi_- : w^T x + b = -1$$

By simple co-ordinate concepts we can get

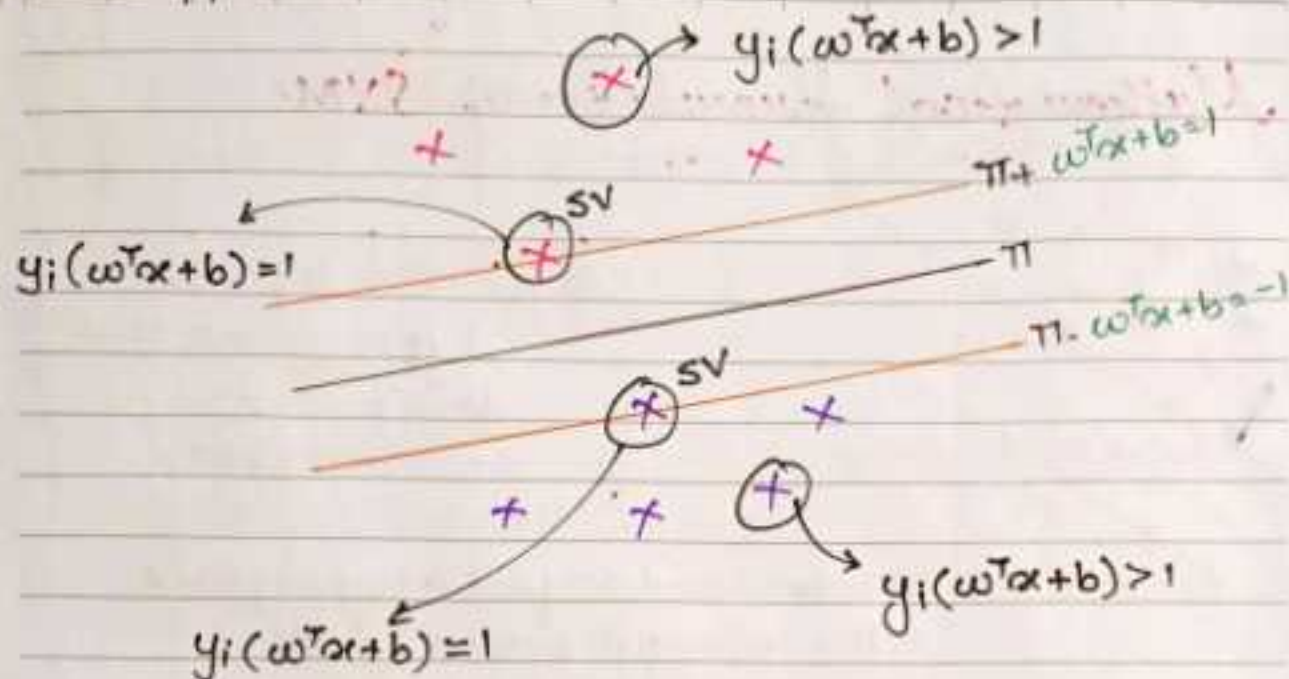
$$d = \frac{2}{\|w\|}$$

So, we have to maximize 'd' with the constraint that all (+ve) points will lie above Π_+ & all (-ve) pts below Π_- .

So, the optimization function becomes

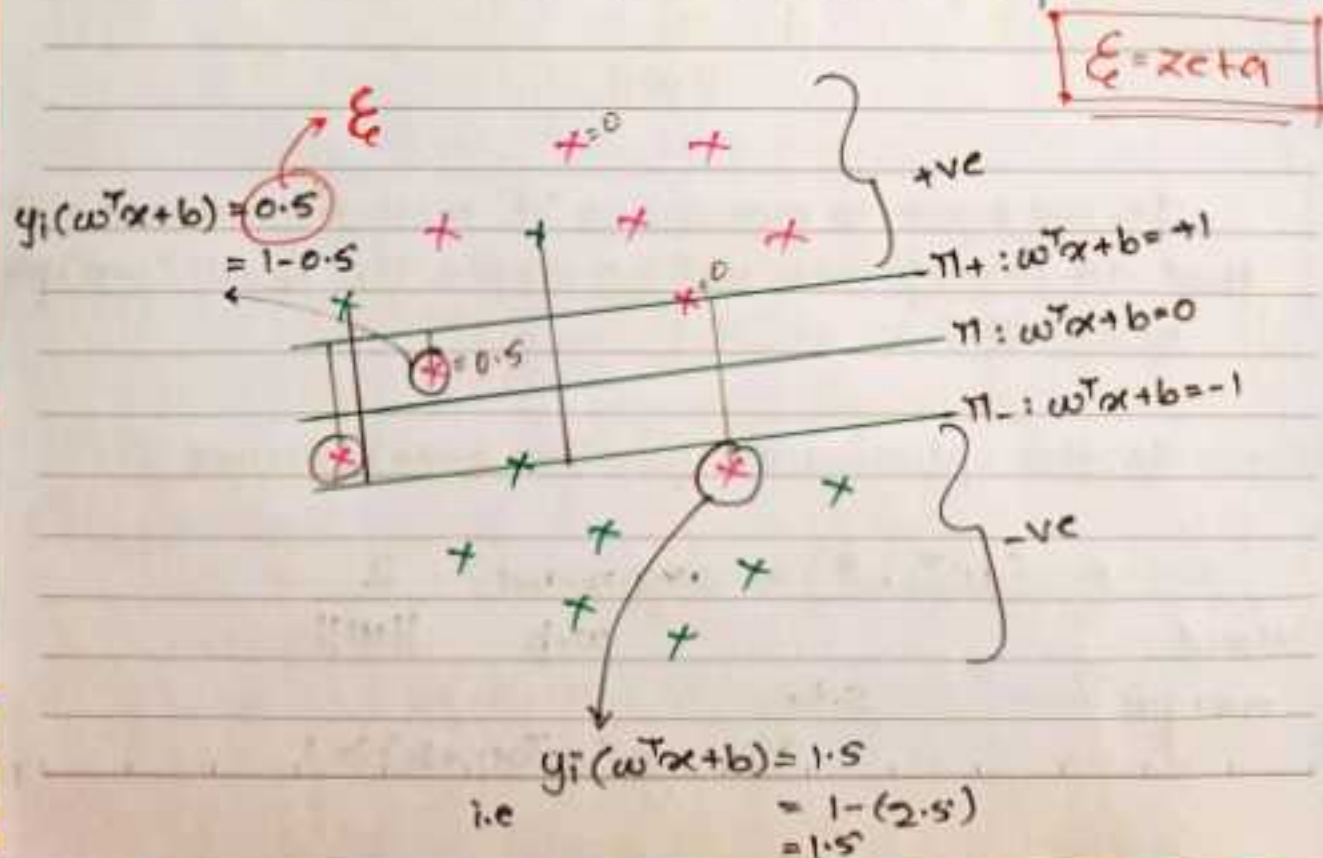
$$\begin{aligned} \text{Hard margin} \left\{ \begin{aligned} (w^*, b^*) &= \arg \max_{w, b} \frac{2}{\|w\|} \\ \text{s.t.} \quad &\forall i, y_i (w^T x_i + b) \geq 1 \end{aligned} \right. \end{aligned}$$

SV: support vectors



But what if we have a dataset that's not linearly separable but almost linearly separable?

→ we formulate an alternative soft-margin constraint that allows some errors to persist.



$$x_i \rightarrow \xi_i$$

$$\xi = \text{zetaeta}$$

$$\xi_i = 0$$

$$\text{If } y_i(w^T x_i + b) \geq 1$$

(i.e. for correctly classified points)

$\xi_i > 0$ and it is equal to some unit of distn away from the correct hyperplane in the incorrect direction.

Now, we want to minimize these errors.

So, our optimization function becomes

$$\underset{w, b}{\text{arg. min}} \quad \frac{\|w\|}{2} + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

margin

avg. distance of misclassified points (hinge loss)

defn of zeta

s.t.

$$y_i(w^T x_i + b) \geq 1 - \xi_i$$

$$\xi_i \geq 0$$

$c \rightarrow$ hyperparameter

$c \uparrow \rightarrow$ tendency to make mistakes \downarrow

\rightarrow overfit (High variance)

$c \downarrow \rightarrow$ underfit (High Bias)

$\forall i$ correctly classified points $\xi_i = 0$

incorrectly classified pts $\xi_i > 0$

maximizing $\frac{2}{\|w\|}$ same as minimizing $\frac{\|w\|}{2}$

$x_i \rightarrow \xi_i$

soft margin of SVM:

ξ_i

$$(w^*, b^*) = \arg \min_{w, b} \frac{\|w\|}{2} + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

s.t.

$$\begin{aligned} y_i (w^T x + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned} \quad \forall i \left. \vphantom{\begin{aligned} y_i (w^T x + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned}} \right\} \text{constraints}$$

away from the correct hyperplane in the incorrect direction.

Now, we want to minimize these errors.

So, our optimization function becomes

$$\arg \min_{w, b} \frac{\|w\|}{2} + c \cdot \frac{1}{n} \sum_{i=1}^n \xi_i$$

margin ← avg. distance of misclassified points (Hinge loss)

s.t.

$$y_i (w^T x_i + b) \geq 1 - \xi_i$$

← defn of zeta

$c \rightarrow$ hyperparameter

$$\xi_i \geq 0$$

$c \uparrow \rightarrow$ tendency to make mistakes \downarrow

\rightarrow overfit (high variance)

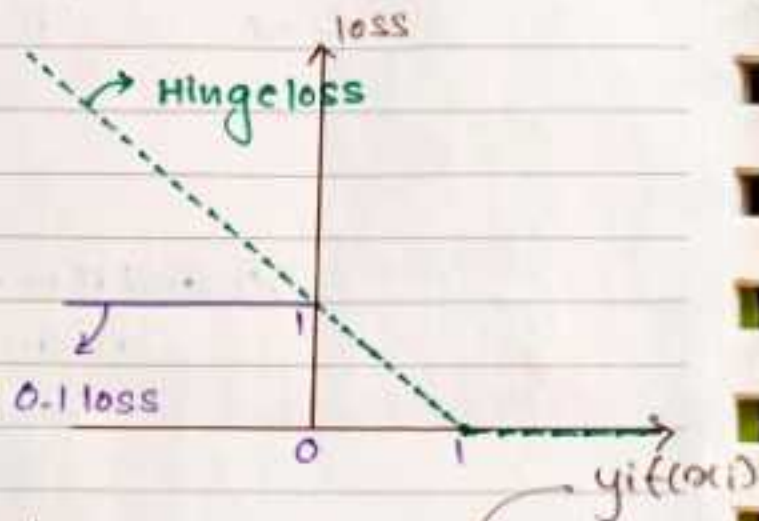
$c \downarrow \rightarrow$ underfit (high Bias)

$\forall i$ correctly classified points $\xi_i = 0$

incorrectly classified pts $\xi_i > 0$

Loss minimization : Hinge Loss

Hinge loss + Regularization
gives us
SVM.



although, Hinge loss also
not differentiable at $\hat{z}_i = 1$,
but since it is continuous unlike
0-1 loss, we can work around
hacks to work with it.

$$\begin{aligned} y_i f(x_i) &= \\ y_i (w^T x + b) \\ &= \hat{z}_i \end{aligned}$$

$\hat{z}_i > 0$: x_i is correctly classified

$\hat{z}_i < 0$: x_i is incorrectly classified.

$$\text{Hinge loss} = \begin{cases} \hat{z}_i \geq 1 ; & \text{Hinge loss} = 0 \\ \hat{z}_i < 1 ; & \text{Hinge loss} = 1 - \hat{z}_i \end{cases}$$

alternatively,

$$\text{Hinge loss} = \max(0, 1 - \hat{z}_i)$$

— Andrew Ng

Lagrange Duality

consider a problem of the following form:

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

we define the Lagrangian to be

$$\mathcal{L}(w, \beta) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Here, β_i 's = Lagrange multipliers.

To find, set \mathcal{L} 's partial derivatives to 0;

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0; \quad \frac{\partial \mathcal{L}}{\partial \beta_i} = 0,$$

and solve for w and β .

consider the following, called primal optimization problem:

$$\min_w \quad f(w)$$

$$\begin{aligned} \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k \\ & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

To solve it, we start by defining
generalized Lagrangian.

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Here, the α_i 's and β_i 's are the Lagrange multipliers.

consider the quantity,

$$\theta_p(w) = \max_{\alpha, \beta; \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta)$$

Here, $p = \text{"primal"}$.

Let some w be given. If w violates any of the primal constraints

(i.e. if either $g_i(w) > 0$ or $h_i(w) \neq 0$ for some i)
then,

$$\begin{aligned} \theta_p(w) &= \max_{\alpha, \beta; \alpha_i \geq 0} f(w) + \sum_{i=1}^k \alpha_i g_i(w) \\ &\quad + \sum_{i=1}^l \beta_i h_i(w) \\ &= \infty \end{aligned}$$

conversely, if the constraints are indeed satisfied for a particular value of w ,

then $\theta_p(w) = f(w)$. Hence

$$\theta_p(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise} \end{cases}$$

Hence, if we consider the minimization problem

$$\min_w \theta_p(w) = \min_w \max_{\alpha, \beta; \alpha \geq 0} \mathcal{L}(w, \alpha, \beta)$$

same as, primal problem.

Optimal value of the objective to be $p^* = \min_w \theta_p(w)$

value of the primal problem.

Now, we define

$$\theta_p(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta)$$

where, $D = \text{"dual"}$

Note: In defⁿ θ_p we were optimizing (maximizing) w w.r.t α, β here we are minimizing w.r.t. w .

• Join me on LinkedIn for the latest updates on Machine Learning: <https://www.linkedin.com/groups/7436898/>

∴ Dual optimization problem:

$$\max_{\alpha, \beta: \alpha_i \geq 0} \Theta_D(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_{\omega} \mathcal{L}(\omega, \alpha, \beta)$$

Optimal value of the dual problem's objective:

$$d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \Theta_D(\omega)$$

How are the primal and the dual problems related?

$$d^* = \max_{\alpha, \beta: \alpha_i \geq 0} \min_{\omega} \mathcal{L}(\omega, \alpha, \beta) \leq \min_{\omega} \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(\omega, \alpha, \beta) = p^*$$

always remember:

"min max min" of a function always being less than or equal to the "min max".

for certain conditions,

$$d^* = p^*$$

i.e. there exists a_i, b_i , so that

$$h_i(w) = a_i^T w + b_i$$

affine \rightarrow linear, except contains
extra intercept term b_i .

Suppose f and the g_i 's are convex, and the h_i 's are affine.

Suppose that the constraints g_i are feasible; means there exists some w so that

$$g_i(w) < 0 \text{ for all } i.$$

under above assumptions, there must exist w^*, α^*, β^* so that w^* is the solution to the primal problem α^*, β^* are the solution to the dual problem, and moreover $p^* = d^* = \mathcal{L}(w^*, \alpha^*, \beta^*)$

moreover, w^*, α^* and β^* satisfy the Karush-Kuhn-Tucker (KKT) conditions

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i=1, \dots, d \quad \text{--- (3)}$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i=1, \dots, l \quad \text{--- (4)}$$

$$\alpha_i^* g_i(w^*) = 0, \quad i=1, \dots, k \quad \text{--- (5)}$$

$$g_i(w^*) \leq 0, \quad i=1, \dots, k \quad \text{--- (6)}$$

$$\alpha^* \geq 0, \quad i=1, \dots, k \quad \text{--- (7)}$$

moreover, if some w^*, α^*, β^* satisfy the KKT conditions, then it is also a solution to the primal and dual problems.

Eqn ⑤, which is called the KKT dual complementarity condition.

Specifically, it implies that if $x_i^* > 0$, then $g_i(w^*) = 0$.

i.e. " $g_i(w) \leq 0$ " constraint is active, meaning it holds with equality rather than with inequality.

This will be key for showing that SVM has only a small number of "support vectors".

Dual ^{Form} of SVM

Following (primal) optimization problem for finding the optimal margin classifier:

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

s.t.

$$y_i (w^T x_i + b) \geq 1, i = 1, \dots, n$$

— (8)

constraints as

$$g_i(w) = -y_i (w^T x_i + b) + 1 \leq 0$$

when we construct Lagrangian for optimization problem we have:

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^n \alpha_i [y_i (w^T x_i + b) - 1]$$

— (9)

Note: there's only " α_i " but no " β_i " Lagrange multipliers, since the problem has only inequality constraints.

To find the dual form, first minimize $\mathcal{L}(w, b, \alpha)$ w.r.t w and b (for fixed α) to get θ_p .

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$\therefore \left[w = \sum_{i=1}^n \alpha_i y_i x_i \right] \quad \text{--- (10)}$$

derivative w.r.t b ,

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i y_i = 0 \quad \text{--- (11)}$$

If we take the defn of w in eqn (10) and put it into the Lagrangian eqn (9), we get

$$\mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j (x_i)^T x_j - b \sum_{i=1}^n \alpha_i y_i$$

but from eqn (11), last term = 0

$$\therefore \mathcal{L}(w, b, \alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j x_i^T x_j$$

∴ with the constraint $\alpha_i \geq 0$, we obtain

dual optimization problem

$$\max_{\alpha} W(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n y_i y_j \alpha_i \alpha_j \langle x_i, x_j \rangle$$

$$\text{s.t. } \alpha_i \geq 0, \quad i=1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

Important Observations:

1. for every α_i , we have an α_i
2. α_i 's only occur in the form of $\alpha_i x_i$

on solving we get

$$f(x_q) = \sum_{i=1}^n \alpha_i y_i x_i^T x_q + b$$

for non-sv, $\alpha_i = 0$

so, for $f(q) \rightarrow$ only the support vectors matters.

3. $\alpha_i > 0$ only for support vectors else 0.

since x_i always occurs only as $x_i^T x_j \rightarrow x_i \cdot x_j$
cosine-sim $\langle x_i, x_j \rangle$

so, basically cosine similarity value are only required to solve the optimization and ultimately get the model.

we can replace cosine-similarity with any other kind of similarity which makes it more powerful.

Generalized dual form-

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{k(x_i, x_j)}_{\text{called Kernel Function.}}$$

(can be any kind of similarity b/w x_i & x_j)

Even during evaluation x_i occurs only as $x_i^T x_q$

$$f(x_q) = \sum_{i=1}^n \alpha_i y_i \underbrace{x_i^T x_q}_{k(x_i, x_q)} + b$$

Kernel trick \rightarrow replacing $x_i^T x_j$ with generalized similarity function $k(x_i, x_j)$.

Quadratic Programming.

The hard margin and soft margin problems are both convex quadratic optimization problems with linear constraints.

Such problems are known as Quadratic Programming (QP) problems

•• Download Machine Learning:
<https://t.me/AIMLDeepThought>

•• Join me on LinkedIn for the latest updates on Machine Learning: <https://www.linkedin.com/groups/7436898/>