

* MACHINE LEARNING NOTES *



What is machine learning ?



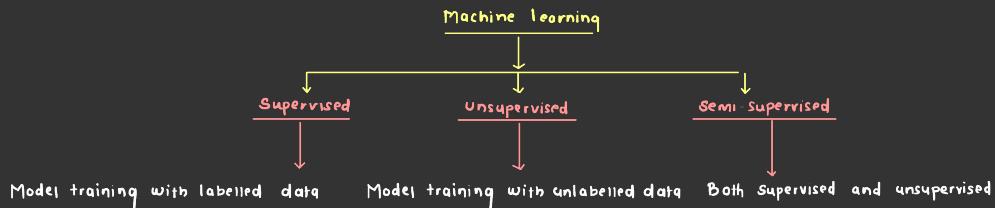
Machine learning is subset of AI.

Statistical model that enabled computer to learn like human

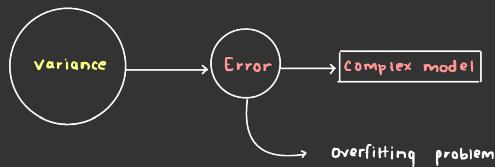
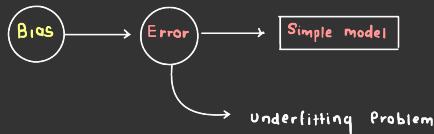
ML rely on data to learn

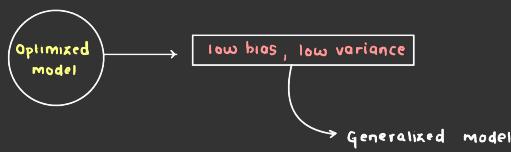
Improve their performance on specific task without being explicitly programmed.

Types of Machine learning:



Bias and variance in machine learning





SIMPLE LINEAR REGRESSION

Simple linear regression statistical analysis that establish linear relationship between x and y .

$x \rightarrow$ independent variable \rightarrow Input

$y \rightarrow$ dependent variable \rightarrow Output

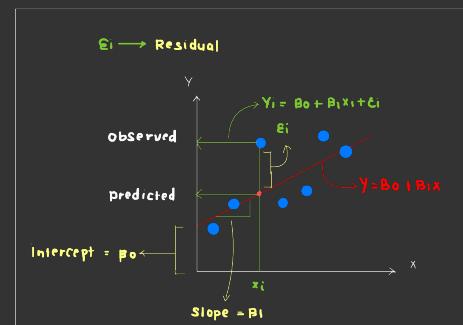
Aim: To find a best fit line using equation $\rightarrow y_i = \beta_0 + \beta_1 x_i$

$\beta_0 \rightarrow$ Constant, intercept

$\beta_1 \rightarrow$ Slope

$x_i \rightarrow$ independent variable

$y_i \rightarrow$ dependent variable



The goal is to find best value for β_0, β_1

Random Error (Residual) $\rightarrow e_i = y_{predicted} - y_i$

\curvearrowleft Error between actual and predicted

where $y_{predicted} = \beta_0 + \beta_1 x_i$

Cost function

$$MSE = \frac{1}{N} \sum_{i=1}^n (y_i - (\beta_1 x_i + \beta_0))^2 \rightarrow \text{Mean Squared Error}$$

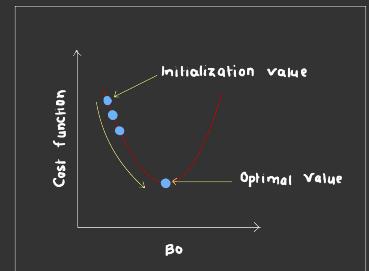
MSE is loss function quantify error than we will update β_0, β_1 parameters using convergence algorithm

Gradient descent → Optimization Algorithm

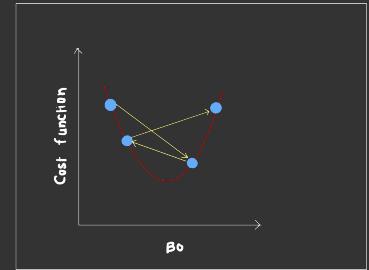
Optimizing → Cost function → Reducing cost function for all datapoints

Updating B_0 and B_1 iteratively until we get optimal solution

The number of steps you are taking to reach optimal value is known as learning rate

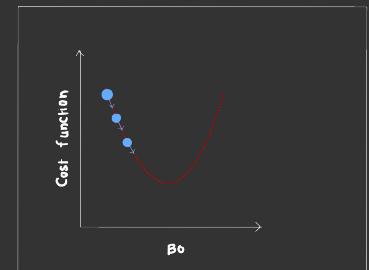


When you choose high of learning rate →



When you choose low of learning rate →

Best value for learning rate → 0.01



To update B_0 and B_1 , we take gradients from cost function. To find gradients, we take partial derivatives for B_0 , B_1 .

$$J = \frac{1}{n} \sum_{i=1}^n (B_0 + B_1 \cdot x_i - y_i)^2$$

$$\frac{\partial J}{\partial B_0} = \frac{1}{n} \sum_{i=1}^n (B_0 + B_1 \cdot x_i - y_i)^2$$

$$\frac{\partial J}{\partial B_1} = \frac{1}{n} \sum_{i=1}^n (B_0 + B_1 \cdot x_i - y_i)^2$$

$$\frac{\partial J}{\partial B_0} = \frac{2}{n} \sum_{i=1}^n (B_0 + B_1 \cdot x_i - y_i) \cdot x_i$$

$$\frac{\partial J}{\partial B_1} = \frac{2}{n} \sum_{i=1}^n (B_0 + B_1 \cdot x_i - y_i) \cdot x_i$$

$$B_0 = B_0 - \alpha \cdot \frac{2}{n} \sum_{i=1}^n (Y_{\text{predicted}} - Y_i)$$

$$B_1 = B_1 - \alpha \cdot \frac{2}{n} \sum_{i=1}^n (Y_{\text{predicted}} - Y_i) \cdot X_i$$

$\alpha \longrightarrow \text{learning rate} \longrightarrow \text{Control convergence Speed}$

Evaluation metrics for linear regression

After training model \longrightarrow Evaluation of model \longrightarrow How accurate model is predicting result.

RMSE \longrightarrow Root mean squared error $= \sqrt{\text{MSE}}$

Lowering RMSE \longrightarrow Better result

RMSE $\longrightarrow 0 \longrightarrow$ A perfect model

Sensitive to outliers

it expresses same unit making it more interpretable.

Coefficient of determination OR R-squared (R²)

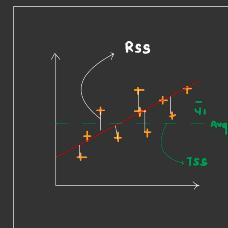
R² always range between $\longrightarrow 0$ to 1

The higher the value of R² \longrightarrow better the model fit data

$$R^2 = 1 - \frac{RSS}{TSS}$$

RSS \longrightarrow Residual sum of squared

TSS \longrightarrow Total sum of squared



$$RSS = \sum_{i=1}^n (Y_{\text{actual}} - Y_{\text{predicted}})^2$$

$\curvearrowright (Y_i - \hat{Y}_i)^2$

$$TSS = \sum_{i=1}^n (Y_{\text{actual}} - \bar{Y}_{\text{predicted}})^2$$

$\curvearrowright (Y_i - \bar{Y})^2$

Problem with R² metric

R^2 tends to increase \uparrow as you add more independent variables. \longrightarrow lead overfitting

\curvearrowleft To solve this we can adjusted R²

Adjusted R² metric

→ Penalizing Overfitting

$$\text{Adjusted } R^2 = 1 - \left[\frac{(1-R^2)(n-k-1)}{(n-k-1)} \right]$$

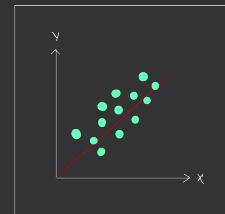
n → no of datapoints

k → no of independent variable

R → R² value

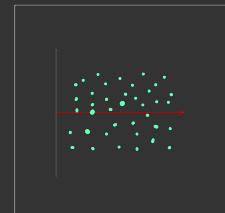
Assumption of linear regression

linear relationship

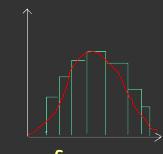


independent of residual

→ The error should not depend on one another.

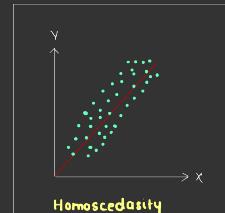


Normal distribution of residual



The equal variance of residual

→ Homoscedasticity



Manual Calculation of Simple linear regression

dataset

X → Height 65 67 70 → inches

Y → Weight 127 160 175 → Pounds

Calculate mean of X and Y

$$\bar{x} = \frac{65 + 67 + 70}{3} = 66.67 \quad \bar{y} = \frac{127 + 160 + 175}{3} = 154.00$$

Calculate the slope

$$\beta_1 = \frac{\sum ((x_i - \bar{x})(y_i - \bar{y}))}{\sum (x_i - \bar{x})^2} \quad \beta_1 = \frac{144.00}{22.33} = 6.45$$

Calculate the intercept

$$\theta_0 = \hat{y} - \beta_1 * \bar{x}$$

$$\theta_0 = 154 - 6.45 * 66.67 = 154 - 430.12 = -276.12$$

Now regression equation

$$\text{Weight} = -276.12 + 6.45 * \text{height}$$

→ We can make prediction using this equation

Calculate the Cost → Error

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

Error for 1st data points → 27759.78

Error for 2nd data points → 178.61

Error for 3rd data points → 4541.98

Hypothesis in linear regression

Whether the coefficient has a statistically meaningful impact on dependent variable
Slope $\rightarrow \beta_1$
 β_1 is significantly different from zero

Whether the independent variable has impact on the dependent variable

Example: $x \rightarrow$ no of hours studied

$y \rightarrow$ Exam's Scores

We want to test whether the no of hours studied is significant predictor of exam scores.
means something important

$H_0 \rightarrow$ no significant relationship between x and y .

or
meaningful

$$\beta_1 = 0$$

$H_1 \rightarrow$ significant relationship between x and y

$$\beta_1 \neq 0$$

Let's say we have collected data from 20 students

Calculate coefficient of hours studied

$$\beta_1 = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} \quad \text{standard error} = \sqrt{\frac{\sum (y_i - \hat{y})^2 / (n-2)}{\sum (x_i - \bar{x})^2}}$$

Let assume $\beta_1 = 0.5$, Standard error = 0.1

Calculate t-stats

$$t\text{-stats} = \frac{(\beta_1 - 0)}{se} = \frac{0.5}{0.1} = 5$$

degree of freedom

$$df = n-2 = 20-2 = 18$$

Now p-value

At $t=5$, $df=18 \rightarrow p\text{-value} = 0.001$

$$\text{Conf. interval} = 95\% = 0.05$$

As we can our p-value is less than conf. interval

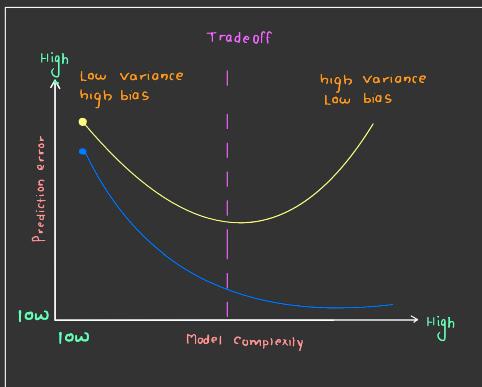
Conclusion

Reject $\rightarrow H_0$: null hypothesis

The coefficient of hours studied β_1 is significantly different from zero.

means x and y have strong relationship.

Bias Variance tradeoff



● → Test data

● → Train data

Tradeoff → Optimized or generalized model

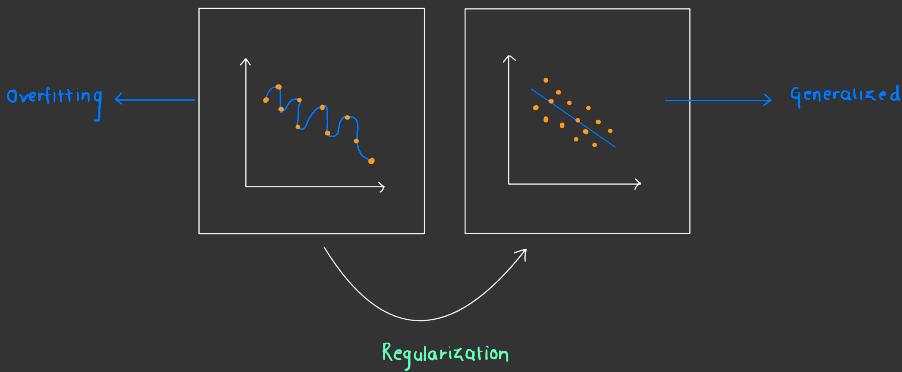
left side area → underfitting

Right side area → overfitting

As we increase model complexity → prediction error
of test data ↑



Regularization → Technique → Prevent → Overfitting → In model.



Ridge Regularization

Add penalty to loss function that is proportional to absolute value of model's coefficient.

encouraging \rightarrow some coefficient \rightarrow to become zero
 \curvearrowright Not exactly zero.

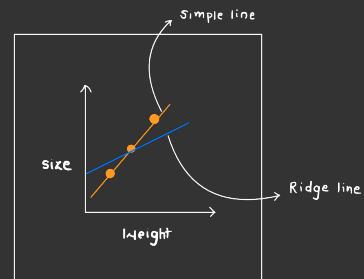
Ridge regularization \rightarrow also known as \rightarrow L2 regularization.

$$\text{Cost function} = \text{loss} + \lambda * \sum_{i=1}^n (\beta_i)^2$$

loss \rightarrow cost function \rightarrow MSE

λ \rightarrow strength of regularization or hyperparameter

β_i \rightarrow slope \rightarrow Model's coefficient.



As λ increases \rightarrow Coefficient leads to shrinkage \rightarrow towards the zero

Ridge \rightarrow introducing \rightarrow bias in model

Lasso Regularization

Add penalty term to linear regression equation

encouraging \rightarrow model to have \rightarrow Smaller coefficient

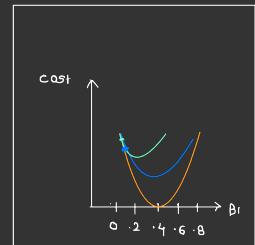
or

even push some of them to become \rightarrow Exactly Zero.

\curvearrowright Some of features will get discarded

$$\text{Cost function} = \text{loss} + \lambda * \sum_{i=1}^n |\beta_i|$$

\curvearrowright We can say lasso can be used as feature selection



λ \rightarrow Responsible to set coefficient value exact zero as well \rightarrow Not all the time but sometime

loss \rightarrow cost function

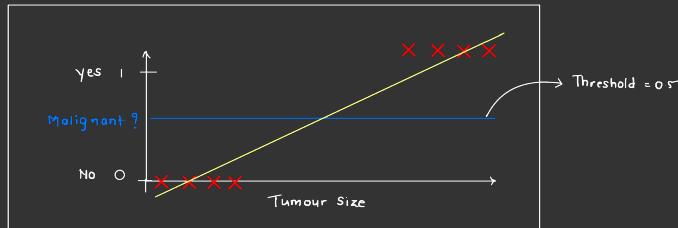
λ \rightarrow strength of regularization

β_i \rightarrow slope

LOGISTIC REGRESSION

logistic Regression is a statistical model used for binary classification
 $\rightarrow (0,1), (\text{yes}, \text{no})$

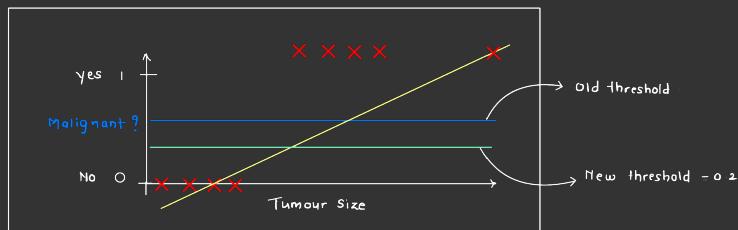
Why logistic regression?



if $y \rightarrow$ greater than $\rightarrow 0.5$ then outcome $\rightarrow 1 \rightarrow$ Malignant tumour

if $y \rightarrow$ less than $\rightarrow 0.5$ then outcome $\rightarrow 0 \rightarrow$ Benign tumour

let suppose in our dataset we have outlier



Now you can see our linear regression is affected by outliers due to which probability may exceeds 1 or go below 0.

To solve this we need logistic regression for binary classification problem.

logistic regression \rightarrow always give output $\rightarrow (0,1)$

logistic function

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}}$$

Equation of best fit line $\rightarrow y = \beta_0 + \beta_1 x$

let say $y \rightarrow$ Probability \rightarrow But this will exceed 1 or go below 0

but we know the range of probability $\rightarrow 0-1$


To overcome this we take "Odd" of p

$$P = \beta_0 + \beta_1 x$$

$$\frac{P}{1-P} = \beta_0 + \beta_1 x$$

 Odd will always post \rightarrow will range between $(0, +\infty)$

Odd is nothing but ratio of $\rightarrow \frac{\text{Prob of success}}{\text{Prob of failure}}$

We know that range is restricted

 problem if we restrict the range then our correlation will decrease

 To control this we can take log of odd

$$\log\left(\frac{P}{1-P}\right) = \beta_0 + \beta_1 x$$

Now we can take exponent on both side \rightarrow To predict outcomes

$$\exp\left[\log\left(\frac{P}{1-P}\right)\right] = \exp(\beta_0 + \beta_1 x)$$

$$e^{\ln\left[\frac{P}{1-P}\right]} = e^{(\beta_0 + \beta_1 x)}$$

$$\frac{P}{1-P} = e^{(\beta_0 + \beta_1 x)}$$

$$P = e^{(\beta_0 + \beta_1 x)} - pe^{(\beta_0 + \beta_1 x)}$$

Take P common

$$P = P \left[\frac{e^{(\beta_0 + \beta_1 x)}}{1 - e^{(\beta_0 + \beta_1 x)}} \right]$$

$$1 = \frac{e^{(\beta_0 + \beta_1 x)}}{P} - e^{(\beta_0 + \beta_1 x)}$$

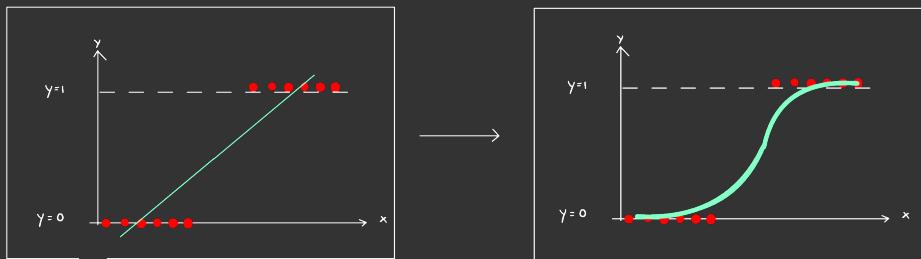
$$P \left[1 + e^{(\beta_0 + \beta_1 x)} \right] = e^{(\beta_0 + \beta_1 x)}$$

$$P = \frac{e^{(\beta_0 + \beta_1 x)}}{1 + e^{(\beta_0 + \beta_1 x)}}$$

now divide by $e^{(\beta_0 + \beta_1 x)}$

$$P = \frac{1}{1 + e^{-(\beta_0 + \beta_1 x)}} \quad \begin{cases} \text{Sigmoid function} \\ \text{Also known as logistic} \end{cases}$$

Now our linear reg become \longrightarrow logistic regression



Cost function

In logistic regression we are not dealing with linearity as you can see above graph.



Problems :

Complicated to reach at \longrightarrow global minima \longrightarrow Because we also getting local minima

To solve this we use different cost function called \rightarrow log loss

log loss \rightarrow derived from the maximum likelihood estimation method

$$\text{log loss} = \frac{1}{N} \sum_{i=1}^N (\gamma_i * \log(\hat{\gamma}_i)) + (1-\gamma_i) * \log(1-\hat{\gamma}_i)$$

What is use of maximum likelihood estimator?

Main aim is to find value of parameters and maximize the likelihood function

estimates the coefficient values

β_0, β_1, \dots

likelihood function \rightarrow In simple word its like Scorecard. It tells you how good or bad your model is explaining real world data

Two Outcomes \rightarrow Success or failure

Bernoulli trials

$$y \sim \text{Ber}(p)$$

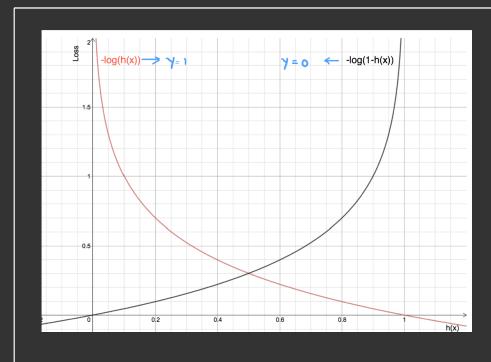
Where the p \rightarrow Sigmoid function

$$L(\beta_0, \beta_1, \beta_2, \dots, \beta_k) = \prod_{i=1}^n p_i^{\gamma_i} * (1-p_i)^{1-\gamma_i}$$

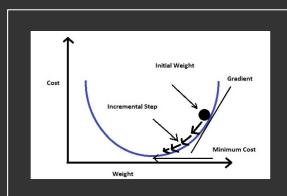
$p_i \rightarrow$ predicted probability

$\gamma_i \rightarrow$ Actual outcomes

$$\log L = \sum [\gamma_i * \log(p_i) + (1-\gamma_i) * \log(1-p_i)]$$
$$p_i = \hat{\gamma}_i \rightarrow \text{Both are same}$$



Gradient descent optimization



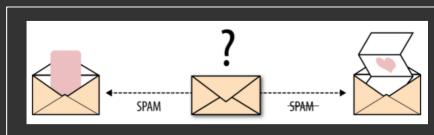
$$\text{Update Rule} \rightarrow \theta_{\text{new}} = \theta_{\text{old}} - \eta \frac{\partial J(\theta)}{\partial \theta_j}$$

$$\frac{d}{dx} \sigma(x) = \frac{d}{dx} \left[\frac{1}{1 + e^{-x}} \right] = \frac{d}{dx} (1 + e^{-x})^{-1}$$

$$= -1 * (1 + e^{-x})^{-2} (-e^{-x})$$

$$\begin{aligned} &= -\frac{-e^{-x}}{(1 + e^{-x})^2} \\ &= \frac{e^{-x}}{(1 + e^{-x})^2} \\ &= \frac{1}{(1 + e^{-x})} \cdot \frac{e^{-x}}{(1 + e^{-x})} \\ &= \frac{1}{(1 + e^{-x})} \cdot \frac{(1 + e^{-x}) - 1}{(1 + e^{-x})} \\ &= \frac{1}{(1 + e^{-x})} \left[\frac{(1 + e^{-x}) - 1}{(1 + e^{-x})} \right] \\ &= \frac{1}{(1 + e^{-x})} \left[1 - \frac{1}{(1 + e^{-x})} \right] \\ \frac{d}{dx} \sigma(x) &= \sigma(x)(1 - \sigma(x)) \end{aligned}$$

Matrices to evaluate classification



Accuracy → Measure the percentage of correctly classified instances.

Example : 200 → emails

140 emails → Correctly predicted as spam → TP

50 emails → Correctly predicted as ham → TN

5 emails → Incorrectly predicted as spam but they are not → FP

5 emails → Incorrectly predicted as ham but they are actually spam → FN

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

may be the not best metric in case of imbalanced dataset.

Confusion metric → Table with combination of predicted and actual value

		Actual	
		Positive	Negative
Predicted	Positive	True Positive	False Positive Type 1 Error
	Negative	False Negative Type 2 Error	True Negative

		$\hat{Y} = 0$ NEGATIVE	$\hat{Y} = 1$ POSITIVE
		$Y = 0$ NOT PREGNANT	$Y = 1$ PREGNANT
TRUE NEGATIVE			
FALSE NEGATIVE			

Precision metrics → how many instances are correctly predicted as positive are actually correct.

use precision → When false positive → higher

$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}} \longrightarrow \frac{\text{No of true post}}{\text{No of predicted positive}}$$

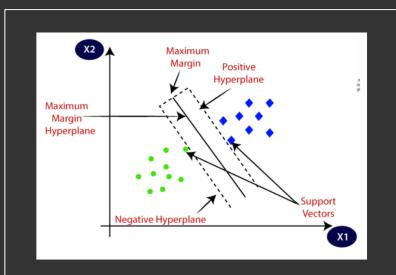
Class imbalance → pos class lesser than neg class

use precision

Support Vector Machine

Support vector machine → supervised machine learning model → used for classification and regression.

Aim → find best hyperplane → separates → data into different class.



Importance Term :

Support vectors → points that are closest to hyperplane

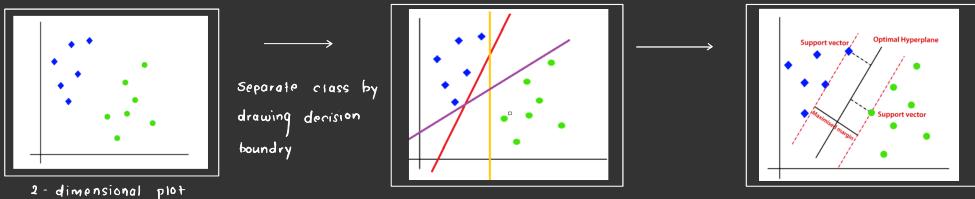
Margin → distance between planes.

Hard margin → aim: perfectly separates the classes

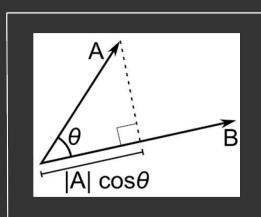
No missclassification in training data
sensitive to outliers. Not always.

Soft margin → Missclassification will be there.

How does SVM work ?



Mathematical Intuition behind SVM



$$A \cdot B = |A| |B| \cos \theta$$

Magnitude of B
Projection of A on B

But in SVM → $|B|$ → Not required

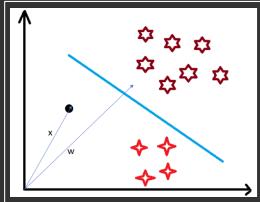
Now $|B|$ → treated as unit vector.

$A \cdot B = |A| |B| \cos \theta$, unit vector of B.

Use of dot product in SVM

Consider \rightarrow point $\rightarrow \vec{x}$

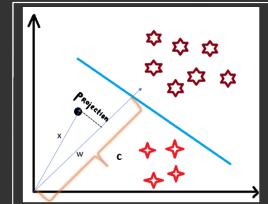
Now we want to know point x lies on left side of plane or right side of plane.



Assume $x \rightarrow$ vector than we make \rightarrow vector w .

$w \perp$ plane

Suppose \rightarrow distance of w from origin to plane $\rightarrow c$

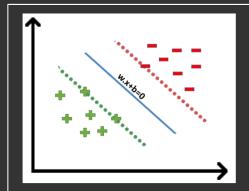


Now $\vec{x} \cdot \vec{w} = c \rightarrow$ point on decision boundary

$\vec{x} \cdot \vec{w} > c \rightarrow$ positive sample

$\vec{x} \cdot \vec{w} < c \rightarrow$ negative sample

Margin in SVM



equation of hyperplane \rightarrow $\vec{w} \cdot \vec{x} + b = 0$

Normal to plane

To classify a point as neg- or post we need decision rule.

$$\vec{x} \cdot \vec{w} - c \geq 0$$

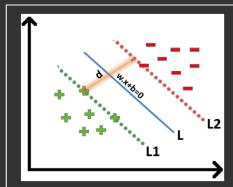
putting $-c$ as b , we get

$$\vec{x} \cdot \vec{w} + b \geq 0$$

hence,

$$y = \begin{cases} +1 & \text{if } \vec{x} \cdot \vec{w} + b \geq 0 \\ -1 & \text{if } \vec{x} \cdot \vec{w} + b < 0 \end{cases} \rightarrow \begin{array}{l} \text{positive} \\ \text{negative} \end{array}$$

Now we need w and b value such that the margin has maximum distance



Let say distance $\rightarrow d$

To calculate distance we need L1 and L2 equation.

Assumptions $\rightarrow \vec{w} \cdot \vec{x} + b = 1 \rightarrow L1$ and $\vec{w} \cdot \vec{x} + b = -1 \rightarrow L2$

with this assumption we will have equal distance from both classes.

Optimization

for all red points $\vec{w} \cdot \vec{x}_i + b \leq -1$
 for all green points $\vec{w} \cdot \vec{x}_i + b \geq 1$

Negative class $\rightarrow y = -1$ Positive class $\rightarrow y = 1$

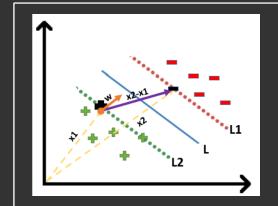
Aim: Maximize the distance between the two planes.

Take Support vector from both classes and distance between these two vectors x_1 and x_2 will be $(x_2 - x_1)$

Finding projection vector on another vector

$$(\vec{x}_2 - \vec{x}_1) \cdot \frac{\vec{w}}{\|\vec{w}\|}$$

$$\frac{\vec{x}_2 \cdot \vec{w} - \vec{x}_1 \cdot \vec{w}}{\|\vec{w}\|} \rightarrow \textcircled{1}$$



for positive point $y = 1$

$$\vec{w} \cdot \vec{x}_1 + b = 1$$

$$\vec{w} \cdot \vec{x}_1 = 1 - b \rightarrow \textcircled{2}$$

for negative point $y = -1$

$$\vec{w} \cdot \vec{x}_1 + b = -1$$

$$\vec{w} \cdot \vec{x}_1 = -b - 1 \rightarrow \textcircled{3}$$

Putting $\textcircled{2}$ and $\textcircled{3}$ in $\textcircled{1}$

$$= \frac{(1-b) - (-b-1)}{\|\vec{w}\|} = \frac{1-b+b+1}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|} = d$$

Hence, the equation which we have to maximize is:

$$\operatorname{argmax}_{(\vec{w}, b^*)} \frac{2}{\|\vec{w}\|} \text{ such that } y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1 \rightarrow \text{Basically this is for Hard margin which means all data are correctly classified}$$

But the above thing is not always possible so we will look into concept of soft margin

Soft margin SVM

We know that $\max(f(x)) \rightarrow$ can write as $\min(-f(x))$

$$\operatorname{argmin}_{(\vec{w}, b^*)} \frac{\|\vec{w}\|}{2} \text{ such that } y_i (\vec{w} \cdot \vec{x}_i + b) \geq 1$$

To make soft margin equation we add 2 more terms $\rightarrow \lambda \alpha, C$

$$\operatorname{argmin}_{(\vec{w}, b^*)} \frac{\|\vec{w}\|}{2} + C \sum_{i=1}^n \xi_i$$

for correctly classified $\rightarrow \xi_i = 0$

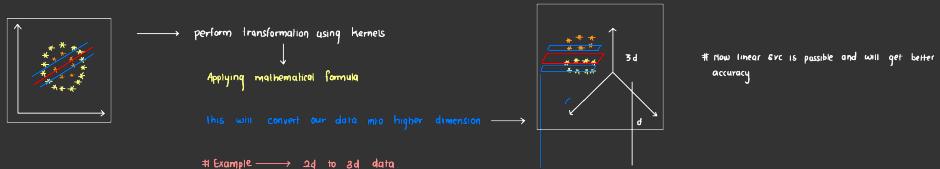
for incorrectly classified \rightarrow distance of that point from its correct hyperplane

Svm kernels

(Importance of kernel) → Transform input data into higher dimensional spaces and making them more easily separable or capturing more complex

relationship among data points

Handles non-linear Separable data



(all dataset)

x	y
1	yes
2	no
3	yes

Let's plot dataset

→ x id Representation

→ Apply Kernel
 $x - x^2$

