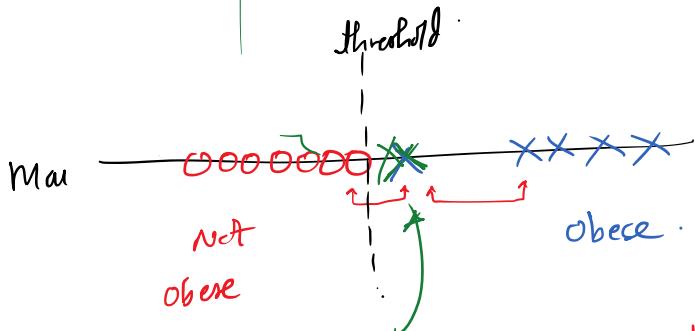


SVM [Support Vector Machine]

Consider \rightarrow Mass of Bandot mice \rightarrow threshold (randomly)

If mass of mice $<$ threshold then declare \rightarrow Not Obese
 " " " " " > threshold " " " " \rightarrow Obese.

But Consider

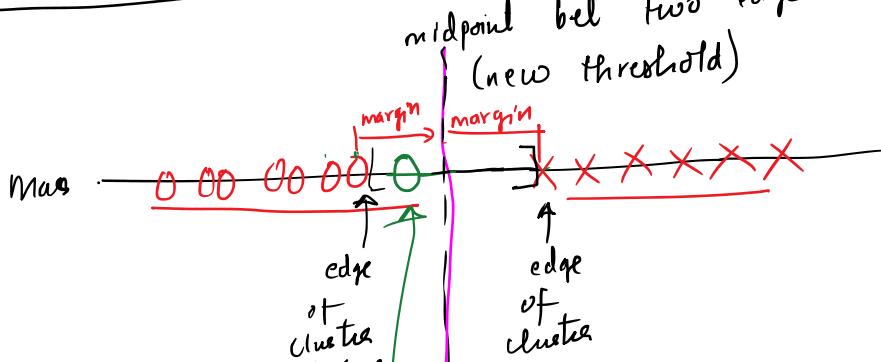


For this observation we will classify it as obese but it is actually closer to observations that are not obese

So the Threshold concept is pretty lame.

We can do better on the original dataset when we focus on the edges of each cluster.

midpoint betⁿ two edge node of each cluster.





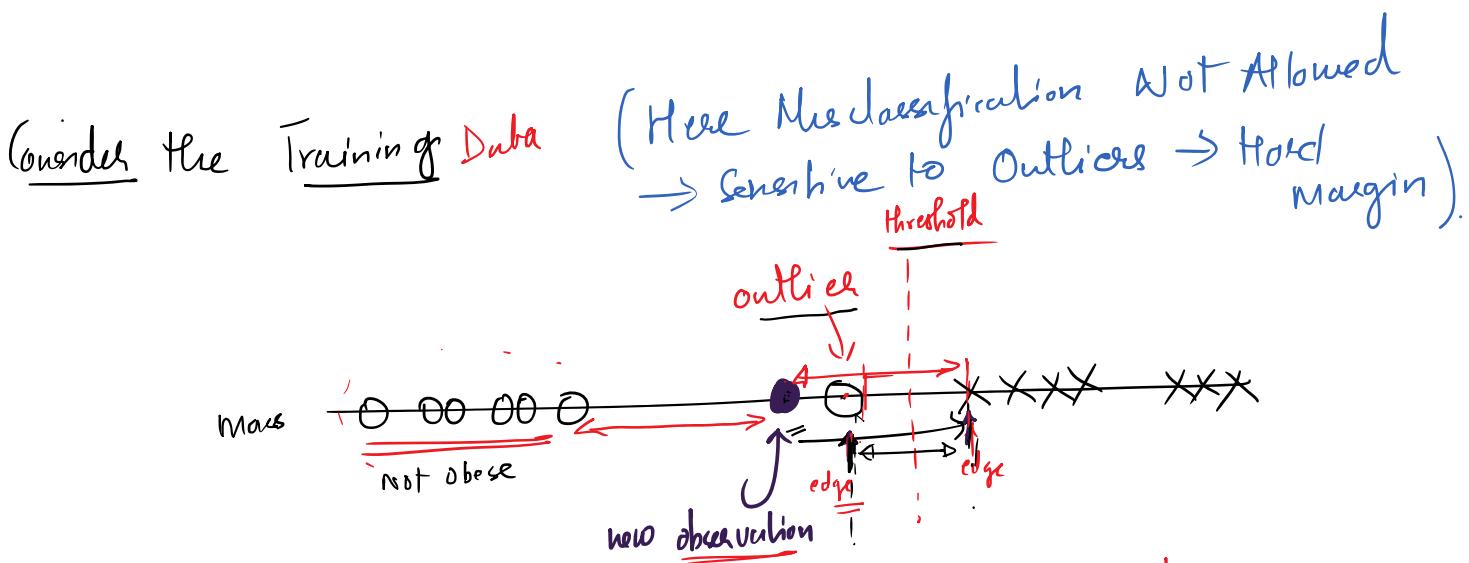
Here if any New observation falls to the left of the mid point then classify it as not obese.

Margin \rightarrow The shortest distance b/w the observations and the threshold is called margin.

Here since threshold is midpoint the margin is same on both sides.

When the threshold is Midway b/w the observations (edge observation of both clusters) then the margin is as large as it can be.

This is known as Maximal Margin Classifier.

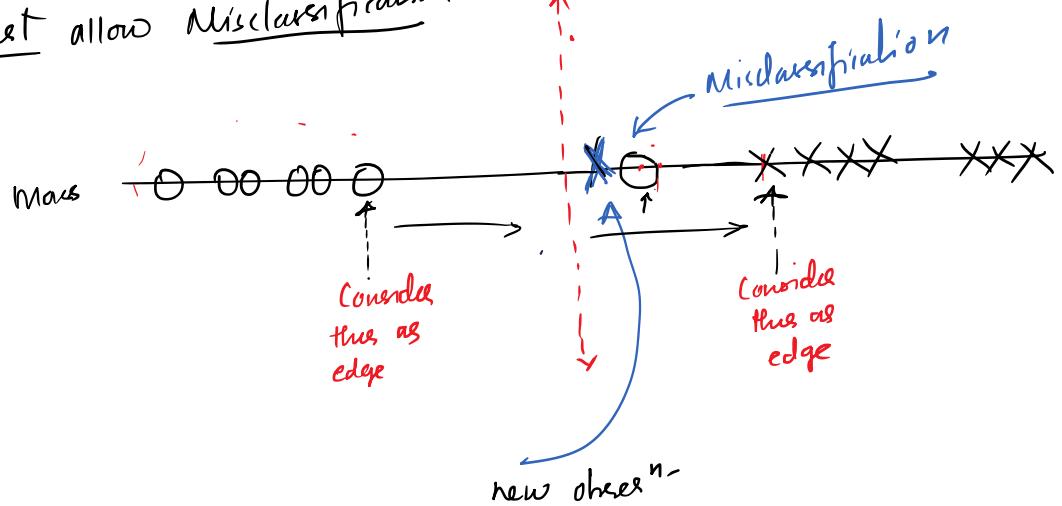


Here for new observation it will be classified as not obese but we can see, it is closer to obese observations.

∴ Maximum Margin Classifier are super sensitive to outliers.

Outliers in Training Data and it makes it pretty ~~un~~

- * We can do better (to make threshold Insensitive to outlier)
For this we must allow Misclassification



Now for the new observation it will be classified as Observe

- * When we pick a threshold that was sensitive to outlier we say it has low bias and it gives error for new observation so we say it has high variance
- * When we select a threshold that allows Misclassification then we say it has high bias but correctly classifies new observation so it has low variance

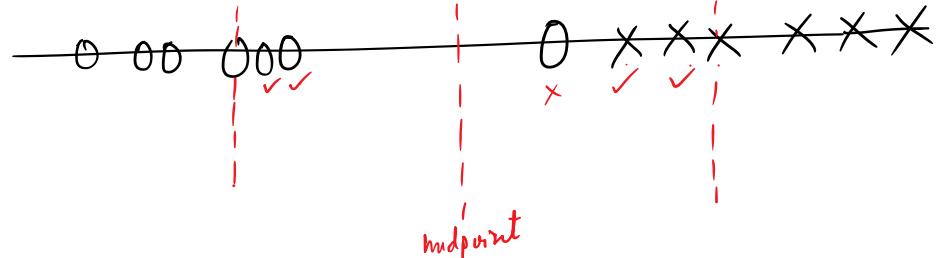
Soft Margin \rightarrow When Misclassification is allowed then the distance b/w threshold and observation is known as Soft Margin

We can't ... normal Misclassification! Allowed, then how to know

We can have many misclassification allowed, then how to know which soft margin is better

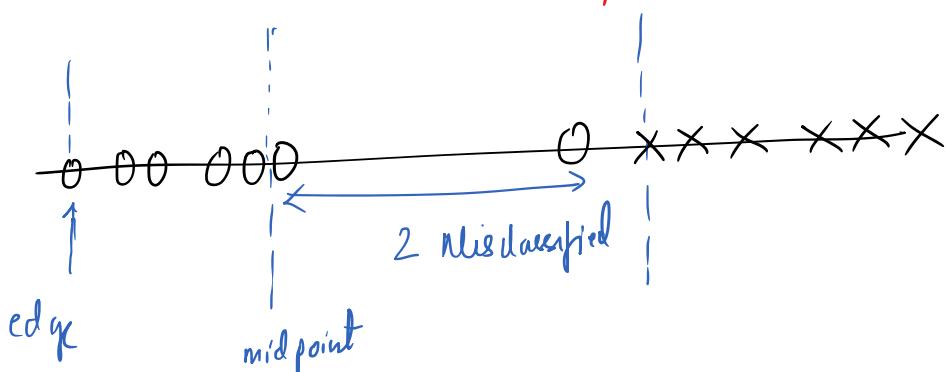
Consider

the soft margin



If are correctly classified
is incorrectly classified

Also For this
soft margin



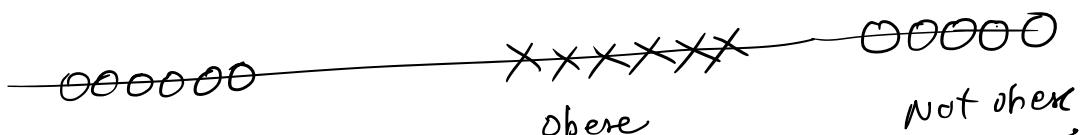
* We will determine which soft margin is better by observing how many misclassification each soft margin

allows.

* Vlm When we use Soft Margin to determine the location of threshold then we are Using "Soft Margin Classifier"

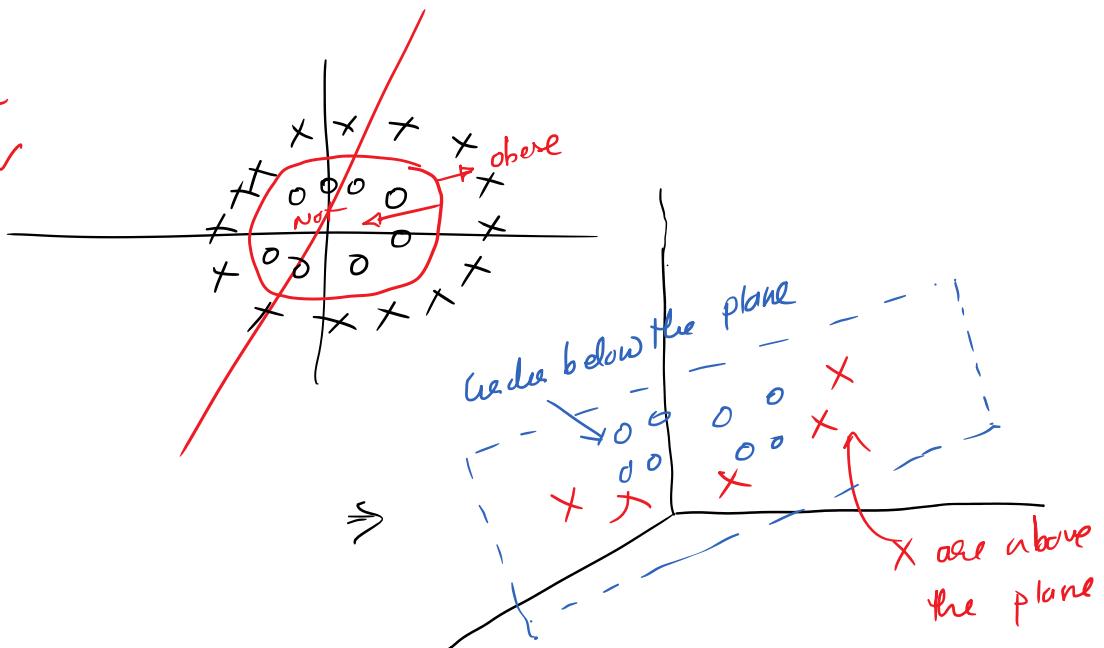
Support Vector Classifier

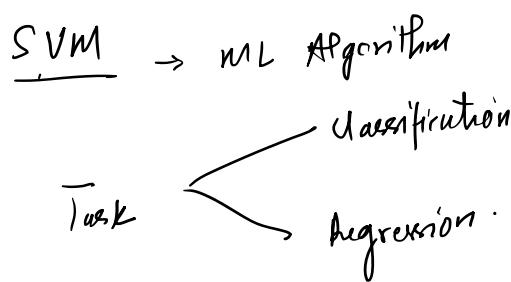
But if the dataset is as.



there Soft Margin will not be efficient as the overlapping is very high.

We cannot have linear Separation
Here

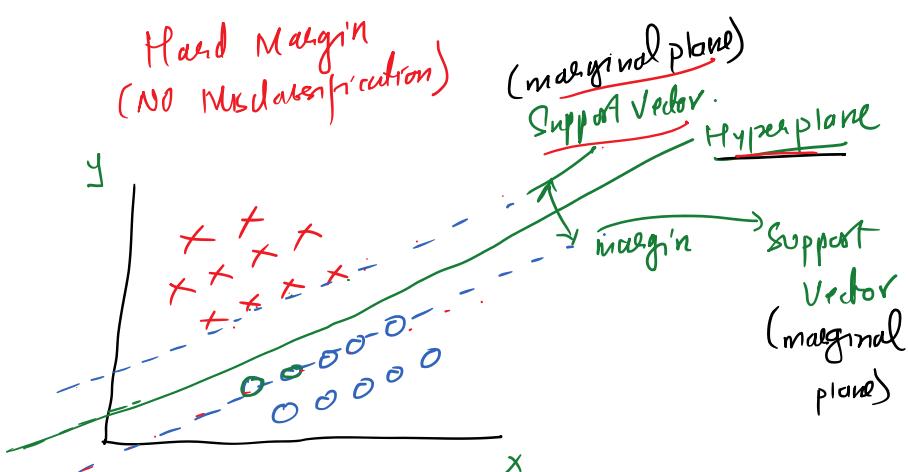




Support Vector Classifier

Consider Binary Classification

2-D Data



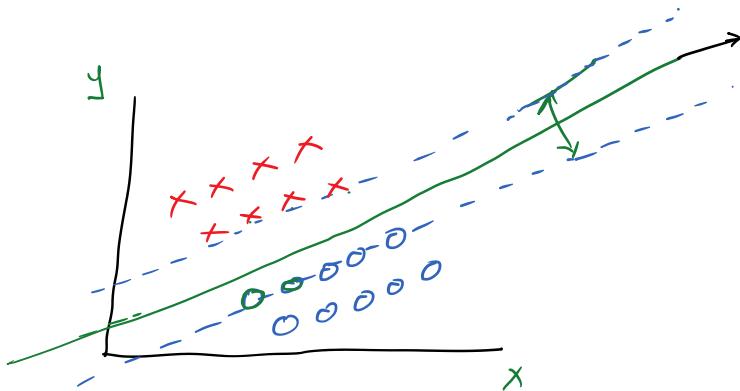
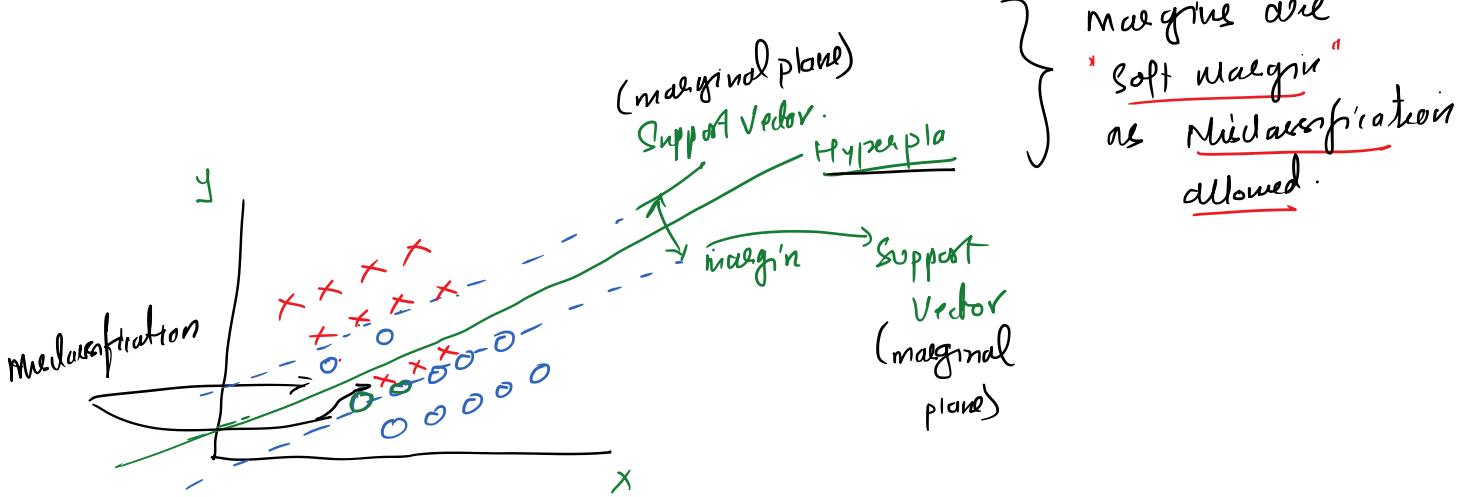
In logistic/linear \Rightarrow Single line (Hyperplane)

SVM \Rightarrow ① Support Vectors (Marginal plane)
② Hyperplane

AIM \rightarrow Create a Hyperplane and Marginal planes such that the margin must be maximum.

* If we get perfect marginal plane and Hyperplane that can linearly separate the Data points then we call it as "Hard Marginal Plane" [No Error / Misclassification]

* If there are Misclassification allowed then the planes (marginal) will be known as "Soft marginal plane"



Note: hyperplane \rightarrow straight line

$$\begin{aligned} \text{Eq} &\rightarrow y = c + mx \\ \text{or} &\quad y = \theta_0 + \theta_1 x \\ \text{or} &\quad y = b + w_1 x \end{aligned} \quad \left. \begin{array}{c} \\ \\ \end{array} \right\} \text{Linear}$$

Multi Linear Regression \rightarrow Many features $x_1, x_2, x_3 \dots$
each with weight w_1, w_2, w_3

$$y = w_1 x_1 + w_2 x_2 + \dots + w_n x_n + b$$

$y = \omega^T x + b$

$w^T \Rightarrow$ weight matrix
 $x =$ feature matrix

Normally Eqⁿ of line $y = mx + c$; - 1

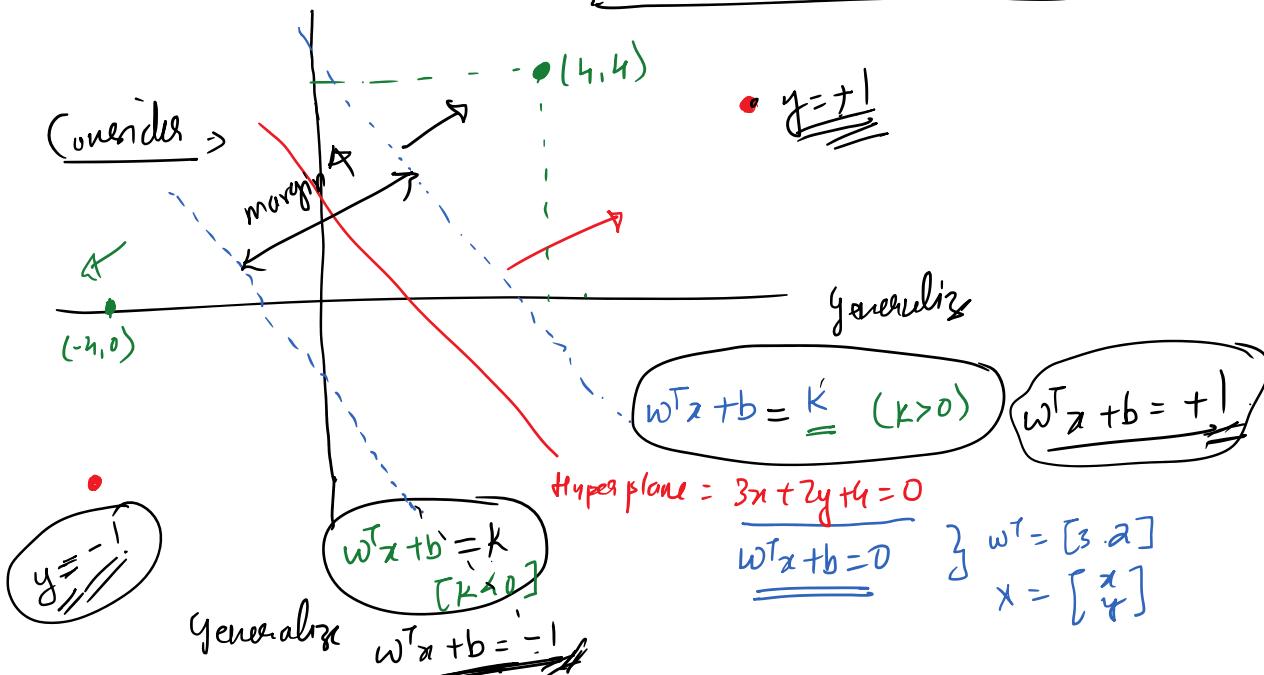
we have $ax + by + c = 0$; - 2

write ② in form of ①

$$\begin{aligned} by &= -ax - c \\ y &= \left(-\frac{a}{b}\right)x - \left(\frac{c}{b}\right) \end{aligned} \quad - 3$$

from ① & ③

$$m = -\frac{a}{b}, \quad c = -\frac{c}{b}$$



(consider the above hyperplane Eqⁿ $3x + 2y + 4 = 0 \rightarrow$ in the form $w^T x + b = 0$)

For the point $(-4, 0) \Rightarrow$ substitute in Eqⁿ

$$3(-4) + 2(0) + 4 = -12 + 0 + 4 = -8 < 0$$

For the point $(4, 4) \Rightarrow$ Substitute in Eqⁿ

$$\begin{aligned} 3(4) + 2(4) + 4 &= 12 + 8 + 4 = 24 > 0 \\ &= \underline{\underline{24}} > 0 \end{aligned}$$

For any point lying above the Eqⁿ $\underline{\underline{w^T x + b = +1}}$ will give +ve value

Similarly any point lying below the Eq $w^T x + b = -1$ will give -ve value.

The hyperplane is represented as $\underline{\underline{w^T x + b = 0}}$

Our Aim to draw two marginal planes (+ve & -ve side) and need to ensure the distance is maximum.

- * The Eqⁿ of marginal plane on +ve side $\underline{\underline{w^T x + b = +1}}$
- * The Eqⁿ of " " " -ve side $\underline{\underline{w^T x + b = -1}}$

- * we want the distance betⁿ them \Rightarrow

lets find diff $\Rightarrow w^T x_1 + b = +1$
 $w^T x_2 + b = -1$
 $\underline{\underline{w^T(x_1 - x_2) = \alpha}} \quad \text{--- (4)}$

$w = \text{coefficient} = \text{slope}$ magnitude $\|w\|$
 vector \vec{w}

We have to focus on Vector.

From w to get $\vec{w} \Rightarrow \frac{w}{\|w\|} \Rightarrow \vec{w}$

From (4)

$$w^T(x_1 - x_2) = \alpha$$

To convert into Vector Divide both side by $\|w\|$

$$\therefore \left(\frac{w^T}{\|w\|} \right) (x_1 - x_2) = \left(\frac{\alpha}{\|w\|} \right)$$

$$\therefore \frac{\vec{w}^T}{\|\vec{w}\|}(x_1 - x_2) = \left(\frac{\alpha}{\|\vec{w}\|} \right)$$

$$\vec{w}^T(x_1 - x_2) = \frac{\alpha}{\|\vec{w}\|}$$

difference/margin

We Need to Maximize Margin \Rightarrow

$$\boxed{\text{Aim} = \underset{(w,b)}{\text{Maximize}} = \frac{\alpha}{\|\vec{w}\|}}$$

By changing w & b

Constraint \Rightarrow

Such that

$$y_i (\text{prediction}) \left\{ \begin{array}{ll} +1 & \text{when } \underline{\vec{w}^T x + b \geq 1} \\ & \text{when point lies outside} \\ & \text{then classification for the} \\ & \text{point } \Rightarrow +1 \\ -1 & \text{when } \underline{\vec{w}^T x + b \leq -1} \\ & \text{i.e. the point lies below the line} \\ & \text{then } \underline{\vec{w}^T x + b = -1} \text{ then } y_i \text{ will be} \\ & \text{classified as } -1 \end{array} \right.$$

+ margin plane

The above constraints are for correctly classified data points

\therefore For all correctly classified data points

Final constraint w.r.t support vector classifier =

$$\boxed{y_i (\text{observed}) * (\vec{w}^T \vec{x}_i + b) \geq 1}$$

$\uparrow \quad \uparrow$
observed predicted

Final \dots

Optimization problem

Final Cost $F^n \rightarrow$

$$\text{Maximize}_{(w,b)} = \frac{\alpha}{\|w\|^2} \quad \text{subjected to constraint} \quad y_i^* \cdot (w^T x_i + b) \geq 1$$

(can be also written as

Since cost F^n we express in terms of Minimize distance bet' margin will increase

$$\text{Minimize}_{(w,b)} = \frac{\alpha}{\|w\|^2} \quad \text{subjected to} \quad y_i^* \cdot (w^T x_i + b) \geq 1$$

But above is for correctly classified points (No Misclassified)

But In real world there is always Misclassification

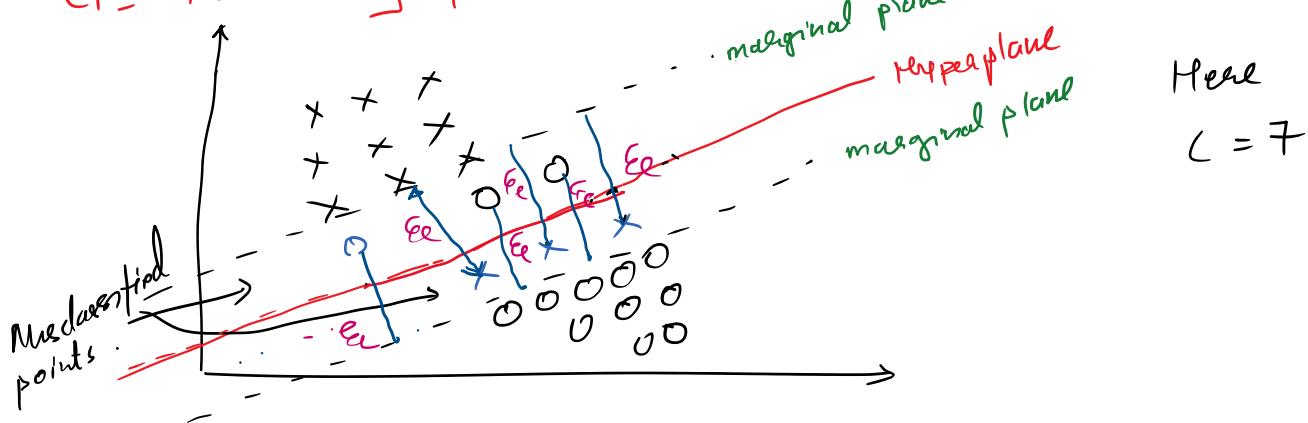
we must take about by using Hyperparameters

∴ final cost F^n (for All Data points (correct/incorrect classified)).

$$\text{Minimize}_{(w,b)} \frac{\alpha}{\|w\|^2} + C_i \sum_{e_i=1}^{C_i} e_i \quad \text{Eta: } C_i = \text{No of misclassified}$$

Hyperparameters.

C_i = How many points we can allow Misclassified.



$C_i \Rightarrow$ It decreases Overfitting problem by allowing

Misclassification

$\sum e_i \Rightarrow$ sum of distance of misclassified points from corresponding margin.

$\sum \epsilon_i$ Measurement \Rightarrow sum of distance of misclassified points Margin

∴ The Cost F^n for SVC is

$$\boxed{\text{Minimize}_{(w, b)} \frac{\|w\|}{\alpha} + C \sum_{i=1}^n |\epsilon_i|}$$

Here Support Vector used for classifⁿ.

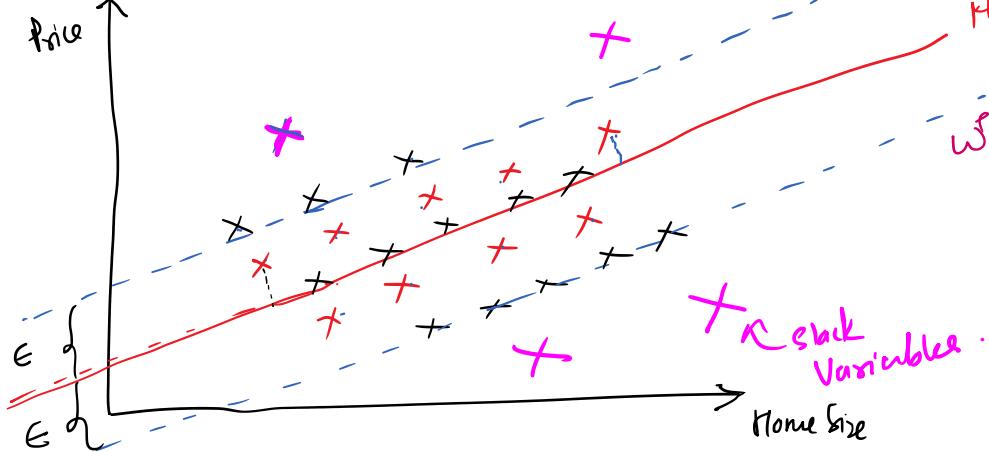
Support Vector for Regression

- * In Linear Regression - we minimize Mean Squared Error (MSE)

Here we are not classifying
we are only interested in knowing
how far the point is from best fit line

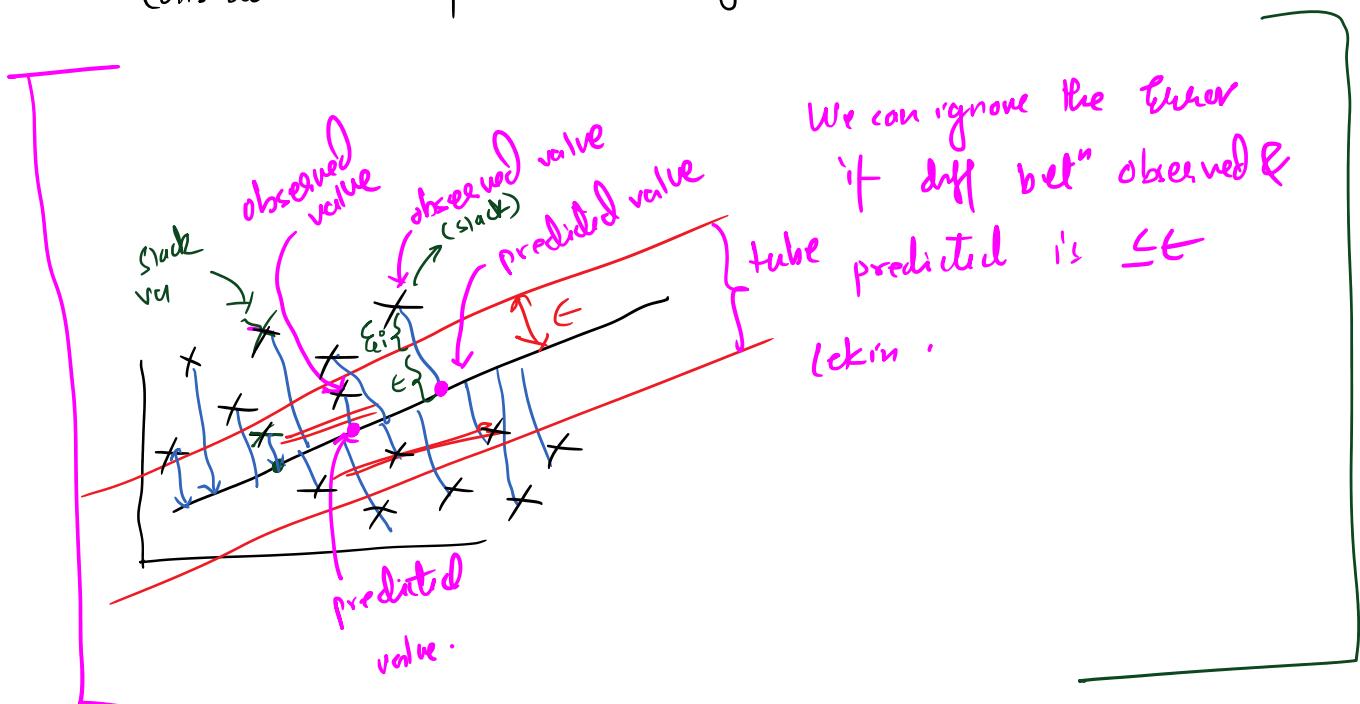
$\text{MSE} = ?$
Best fit line that gives Minimum MSE

In Support Vector Regression



Here we are creating ϵ Insensitive tube (allowed margin of error)

All the points inside ϵ Insensitive tube are not considered for calculating error



For a point x_i^0 within ϵ Insensitive tube

$$|y_i^0 - w^T x_i^0| \leq \epsilon$$

For a point x_i

$$|y_i - w^T x_i| \leq \epsilon$$

y_i = observed value

$w^T x_i$ = predicted value

Constraint if $|y_i - w^T x_i| \leq \epsilon$ then we can say the prediction is good & we will not consider the difference in error calculation.

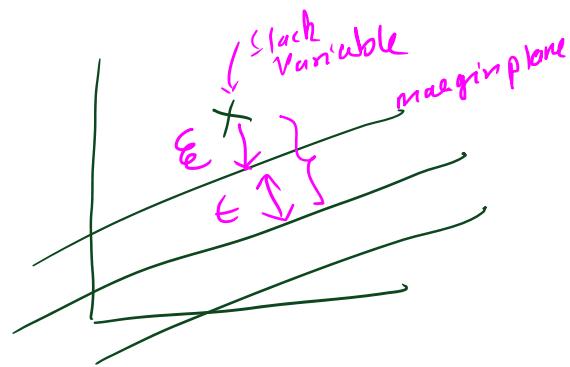
There are points outside marginal plane (Slack Variables)

\therefore we need to calculate distance of slack variables from hyperplane (observed value y_i) (predicted value)

$$= \epsilon + |e_i|$$

\uparrow
distance betⁿ
marginal plane
& Hyperplane

\uparrow
distance betⁿ
marginal plane &
slack variable



$$\text{Cost } J^n = \underset{(w,b)}{\text{Minimize}} \frac{\|w\|^2}{2} + C_i \sum_{i=1}^n |e_i|$$

Hinge Loss

Constraint: $|y_i - w^T x_i| \leq \underline{\epsilon + |e_i|}$

i.e. distance betⁿ observed & predicted value must be

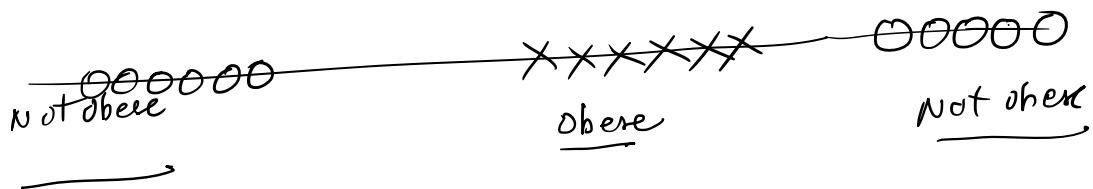
$$\leq \underline{\epsilon + |e_i|}$$

- * As compared to linear regression, the sum is allowing some margin (ϵ) for error or buffer

Margin (ϵ) for error or buffer

- * Here basically the error is distance betⁿ slack variable
(variable outside the tube) and the tube
 ϵ

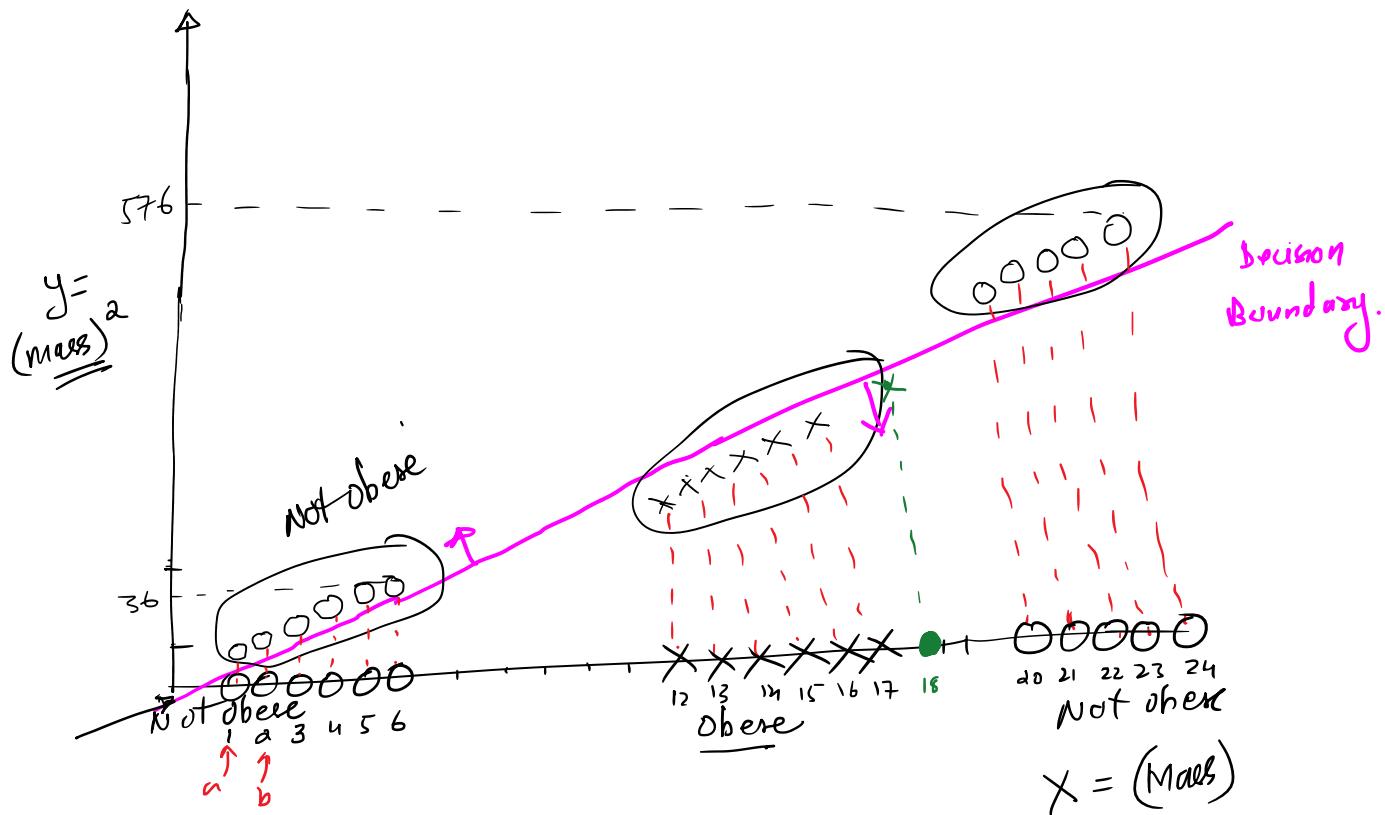
SVM (Support Vector Machine) →
If the Dataset is al.



there Soft Margin will not be efficient as the overlapping is very high.

* Whenever we have lot of Misclassification (Overlapping) then Max Margin Classifier and Support Vector classifier cannot be used.

In such case we will use Support Vector Machine.



Here $X = \text{Mass}$
 $y = (\text{Mass})^2$

from above By increasing dimension of data from 1D to 2D by using Square (Hole) function we can linearly separate the data points.

We can also use it to predict for new observation.

→ For New observation take Square of observation.

SVM (Main Idea)

- * Start data with low Dimension (given)
- * Now Move data to higher Dimension
- * Find SV classifier to classify the data into diff classes.

Question Arised → why Square why not Cube?

i.e How to decide, how to Transform the data.

[SVM Uses Kernel Function to Systematically Find SV Classifier in Higher Dimension]

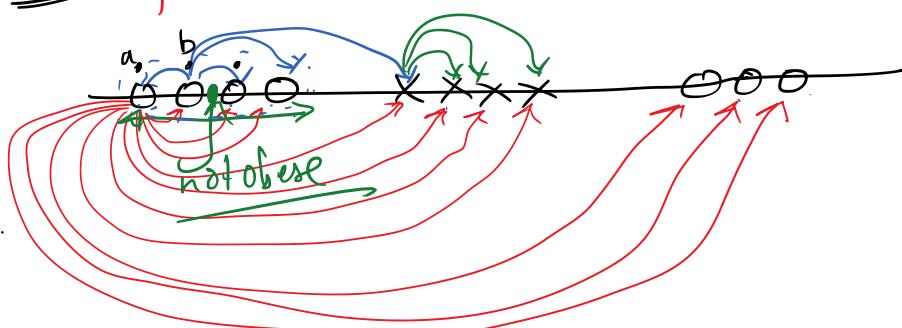
* Vizuf

- * Kernel function actually never does any transformation in Higher Dimension.
 - * Instead of calculate Relation b/w observations and visualizes them in Higher Dimension. SVC uses this Visualizn for classification.
 - * This is Known as Kernel Trick
- * Let us Use Polynomial Kernel that has parameter 'd'
ii. In Low degree of Polynomial.

↗ Let us use memory.

Stands for degree of Polynomial.

Stands for degree of polynomial.
If $d=1$ → the Polynomial Kernel computes relationship betⁿ each pair of observation in 1D



| Here All points with All other points relation is used.

Here d = best found by cross Validation.

* We have Many Kernel f^n

Mast commonly used are →

▷ Polynomial Kernel ↗

- ① Polynomial Kernels
- ② Radial Basis F^n

Polynomial Kernel \Rightarrow It calculates higher dimensional Relationship.

The kernel that transforms 1D to 2D is Polynomial Kernel

It may look like $(a * b + r)^d$

* a and b are two diff observations in dataset
(any two points in dataset).

* r = coefficient of Polynomial.

* d = degree of Polynomial.

[Here we use SVM with Polynomial Kernel to compute relationship betⁿ observations in higher dimension and then find good classification.]

Find Relationship betⁿ observation in higher dimension.

Let $r = 1/2$ $d = 2$ } value of r and d is determined by cross validation.

$$\begin{aligned} (a * b + 1/2)^2 &= (a * b + 1/2) \cdot (a * b + 1/2) \\ &= a^2 b^2 + \frac{1}{2} ab + \frac{1}{2} ab + 1/4. \end{aligned}$$

This polynomial is equal to this dot product.

First observation x-axis values

Second observation x-axis value

higher dimension visual y-axis value.

$$\begin{aligned} &= a^2 b^2 + ab + 1/4. \\ &= ((a, a^2, 1/2)) \cdot ((b, b^2, 1/2)) \} \end{aligned}$$

Dot product is sum of 1st term multiplied, 2nd term multiplied and so on]

* The Dot product gives us high dimensional coordinate for Data.

... ... more

Here γ_2 is third axis but values are same so can ignore

Here $(a+b\gamma)^d$ is used to get higher dimension relⁿ of two data points $a \& b$.

Note
 Since $= (a+b+\gamma_2)^2 = (a, a^2, \gamma_2) \cdot (b, b^2, \gamma_2)$
 \Rightarrow we can observe that all we need to do calculate the higher dimensional relationship is to calculate dot products betⁿ each pair of points.

When
 $r=1, d=2$

$$(a+b+\gamma)^d = (a+b+1)^2$$

$$= a^2b^2 + 2ab + 1$$

$$= 2ab + a^2b^2 + 1$$

$$= (\sqrt{2}a, a^2, 1) \cdot (\sqrt{2}b, b^2, 1)$$

Now $a=1$ $b=14$
 for these two points, we want to find high dimensional relationship betⁿ these two points

Now $(a+b+\gamma_2)^2 = (1+14+\gamma_2)^2 = 126 \cdot 5^2 = \underline{16 \text{ or } 2.25}$

is the one^d-D

(bcz we have to do this for all data points pair)

relationship we need to solve for SV classifier even though we actually did not transform data in d-Dimension

* Once we decide value of r and d using Cross Validation we can find higher dimensional relationship betⁿ every pair of data points.

$v \cdot v$ \vdash v
pair of data points.

Radial Kernel \rightarrow Works in infinite Dimension.

Question \Rightarrow What are the parameters?

\hookrightarrow How it calculates the relationship in Higher Dimension.

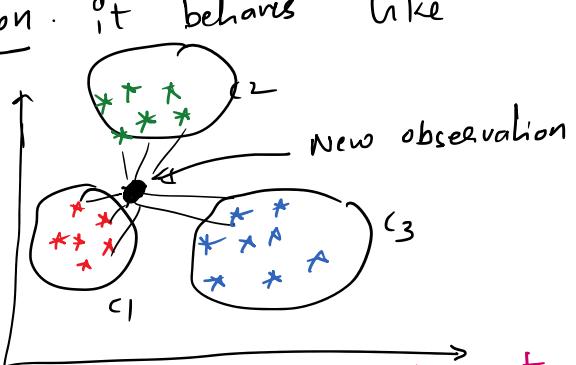
\hookrightarrow It visualizes relationship of observations in infinite Dimension.

- * One way to deal with overlapping data is to use SVM with Radial Kernel.

$$e^{-\gamma(a-b)^2} \quad \left. \right\} \text{Radial Basis Function (RBF)}$$

\rightarrow Since RBF finds SV classifier in Infinite Dimension it is not possible to visualize what it does.

\rightarrow But when applied to New observation it behaves like Weighted Nearest Neighbour Model



* Radial Kernel checks influence of the Nearest Neighbours on New observations and accordingly makes classification.

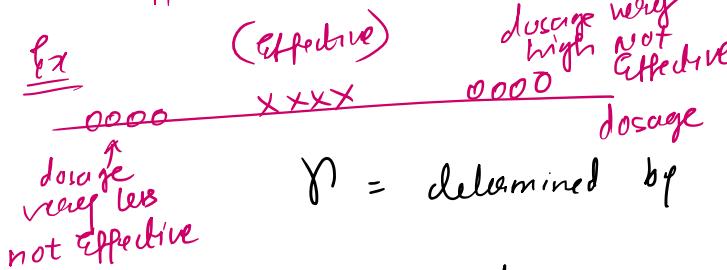
[\star closest observation (Nearest Neighbour will have lot of influence on how we classify new observations as compared to far off observations]

* Let's talk about how radial kernel determines how much influence each observation in Training Dataset has on classifying new observations.

$e^{-\gamma(a-b)^2}$ Here a and b are two observations.
 \downarrow $\Rightarrow (a-b)^2 \Rightarrow$ diff betⁿ the measurement is squared giving us the squared distance betⁿ the two observations.

[let say a and b are ...]

Let say a and b are different data points.
(different dosages)

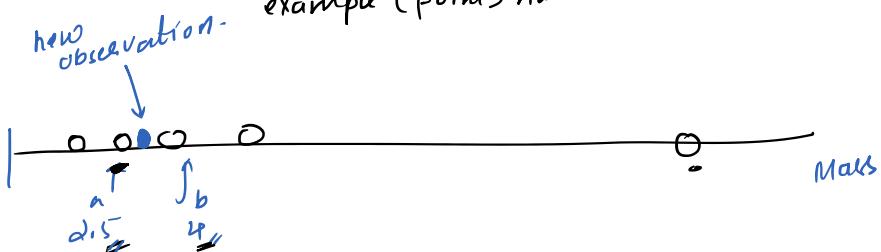


* So the amount of influence that one observation has on another is function of $(a-b)^2$. [Squared distance].

γ = determined by cross validation [γ scales the distance hence influence].

$$= \frac{1}{\text{no of features} * 5^2}$$

= scalar that defines how much influence a single training example (point) has.



The observation $a = 2.5$ and $b = 4$ are closest to the new observation
[These observations are relatively near]

Radial Kernel uses this classification for new observations

Let $\gamma = 1$ γ scales the distance
Here low scale, low distance, overall high influence

$$\bar{e}^{-(2.5-4)^2} = 0.11$$

- (1) When γ is \uparrow

it scales (multiplies) distance betw a & b .
So effectively effects the amount of influence the two points have on each other.

$$\text{let } \gamma = 2$$

$$\bar{e}^{-(2.5-4)^2} = 0.01$$

[0.11 is high dimensional relationship betw

observations]

If two observations are $a = 2.5$ $b = 16$

Here the observations are relatively far from each other.

$$\text{let } \gamma = 1 \quad \bar{e}^{-(2.5-16)^2} = \bar{e}^{-182.25}$$

= close to zero

[This relationship is close to zero as the points are far from each other.]

From (1) & (2)

if points are close enough then we have high influence.

and if points are far off then we have low interaction.
 Just like Polynomial Kernel when we plug the values in
 Radial Kernel we get High Dimensional Relationship.

Note →

To understand How Radial Basis Kernel works in Infinite Dimensions:

Let us take Polynomial Kernel →

$$(a \cdot b + r)^d = \begin{cases} d=0 & \Rightarrow (a \cdot b)^0 = (a)^0 \cdot (b)^0 \\ d=2 & \Rightarrow \text{So new coordinate is only the square of original measurement on original axis} \\ & [\text{No new axis generated}]. \end{cases}$$

when $\begin{matrix} a=d \cdot 1 \\ b=r \end{matrix} \Rightarrow \begin{matrix} a^2 = b \cdot 1^2 \\ b^2 = 1b \end{matrix}$

when $r=0 \quad d=3 \Rightarrow a^3 \cdot b^3$

$r=0 \quad d=1 \quad a \cdot b$

$r=0 \quad d=2 \quad a^2 \cdot b^2$

when $r=0$ leaves data in same 1 dimensional irrespective of values of d .

To add $d=1, d=2 \Rightarrow a' \cdot b' + a^2 \cdot b^2 \Rightarrow \underline{(a \cdot a^2)} \cdot \underline{(b \cdot b^2)}$

This is 2-D Transformation.

More actually no transformation was done.

* Just solved the dot product to get high Dimensional

* just some

relationship.

lets do for ∞ dimension \Rightarrow

$$a^1 b^1 + a^2 b^2 + \dots + a^\infty b^\infty \Rightarrow (a^1, a^2, a^3, \dots, a^\infty) \cdot (b^1, b^2, \dots, b^\infty)$$

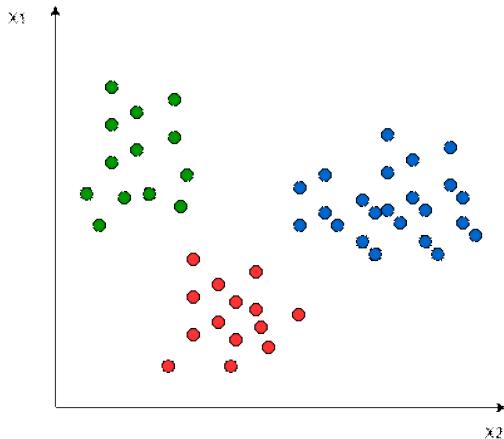
This gives us Dot product with
Coordinate for infinite no of Dimensions.

*
This is what Radial Kernel Does

Multiclass Classification Using SVM

- In its most simple type, SVM doesn't support multiclass classification natively.
- It supports binary classification and separating data points into two classes.
- For multiclass classification, the same principle is utilized after breaking down the multiclassification problem into multiple binary classification problems.
- The idea is to map data points to high dimensional space to gain mutual linear separation between every two classes.
- This is called a [One-to-One approach](#), which breaks down the multiclass problem into multiple binary classification problems. A binary classifier per each pair of classes.
- Another approach one can use is [One-to-Rest](#). In that approach, the breakdown is set to a binary classifier per each class.
- A single SVM does binary classification and can differentiate between two classes.
- So that, according to the two breakdown approaches, to classify data points from classes data set:
- In the *One-to-Rest* approach, the classifier can use SVMs. Each SVM would predict membership in one of the classes.
- In the *One-to-One* approach, the classifier can use SVMs.

Let's take an example of 3 classes classification problem; green, red, and blue, as the following image:

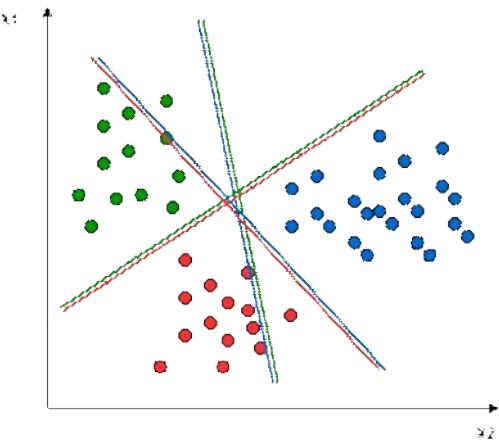


Applying the two approaches to this data set results in the followings:

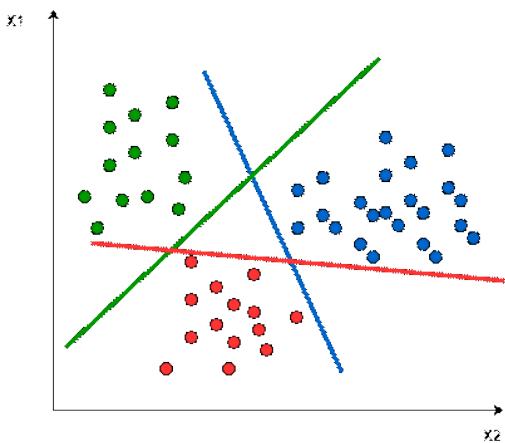
In the *One-to-One* approach, we need a hyperplane to separate between every two classes, neglecting the points of the third class.

This means the separation takes into account only the points of the two classes in the current split.

For example, the red-blue line tries to maximize the separation only between blue and red points. It has nothing to do with green points:



In the *One-to-Rest* approach, we need a hyperplane to separate between a class and all others at once. This means the separation takes all points into account, dividing them into two groups; a group for the class points and a group for all other points. For example, the green line tries to maximize the separation between green points and all other points at once:



One of the most common real-world problems for multiclass classification using SVM is text classification. For example, classifying news articles, tweets, or scientific papers.