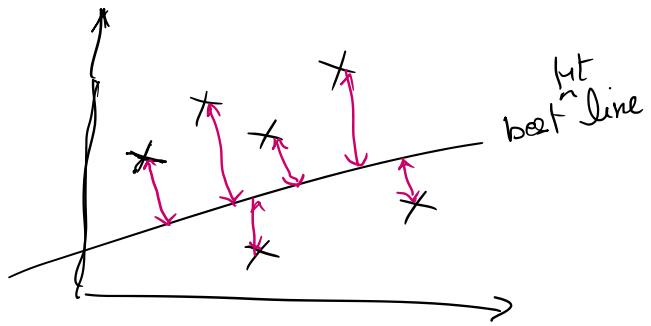


Support Vector For Regression \Rightarrow

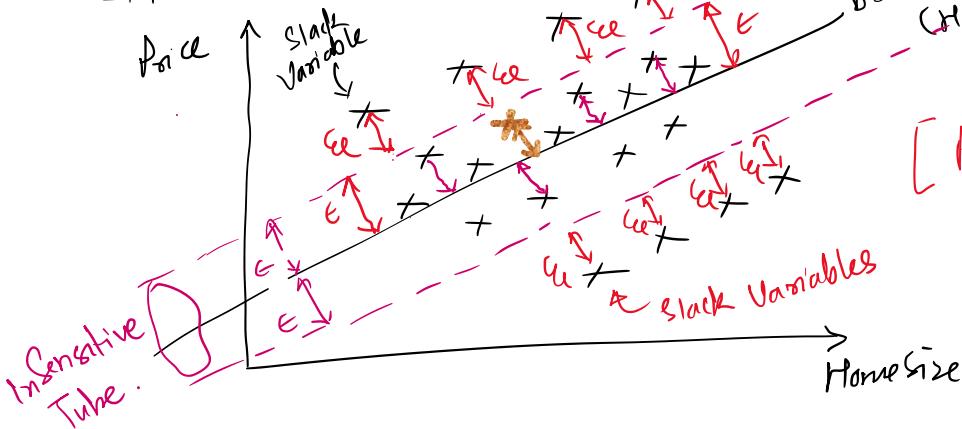
- * In Linear Regression

Objective \rightarrow Minimize MSE
[Mean Squared Error]



Best Fit line $\Rightarrow g_t$ is the line that gives Minimum MSE

Support Vector Regression \Rightarrow



$w^T x + \epsilon$ (margin plane)
 $w^T x$ (margin plane)
best fit line (Hyperplane).
 $w^T x - \epsilon$.

[Data Points Inside ϵ Insensitive Tube \rightarrow Support Vector]

\rightarrow Here we are calculating ϵ (Insensitive Tube) [Allowed Margin of Error]

- * We can ignore the error if difference between observed and predicted value is $\leq \epsilon$

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

The data points outside the ϵ insensitive tube are known as Slack Variables.

- * For a point x_i within Insensitive (ϵ) tube ... normal value

- * For a point x_i within Insensitive (ϵ) tube

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

y_i^* = observed value
 x_i^* = predicted value

constraint

$$\text{if } |y_i^* - w^T x_i^*| \leq \epsilon$$

We can say that prediction is good and we will not consider the difference in error calculation.

- * For Slack Variable [points outside ϵ Insensitive tube]

we need to calculate distance of slack variable from Hyperplane

$$= \text{distance betn slack variable} + \epsilon$$

& Marginal plane

$$= \epsilon_e + \epsilon$$

objective \Rightarrow

$$\text{Cost } F^n \Rightarrow \text{Minimize}_{(w, b)}$$

$$\frac{\|w\|^2}{2} + C_i^* \sum_{e_i=1}^{l_i} \epsilon_e^*$$

Constraints

$$|y_i^* - w^T x_i^*| \leq \epsilon + \epsilon_e^*$$

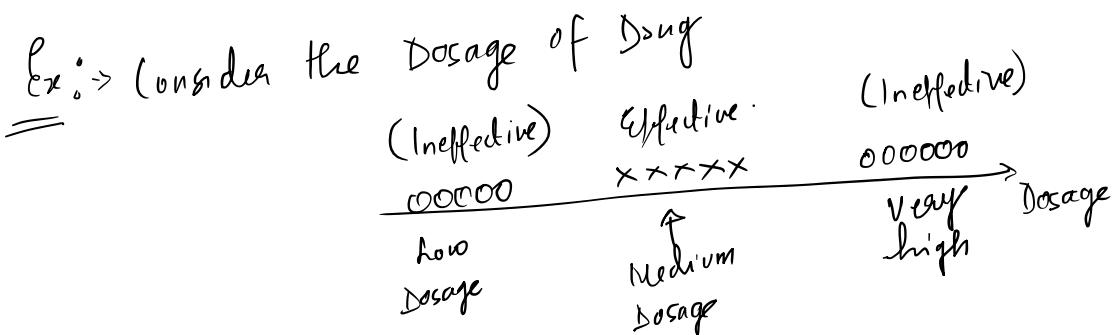
Hinge Loss

Hyperparameter

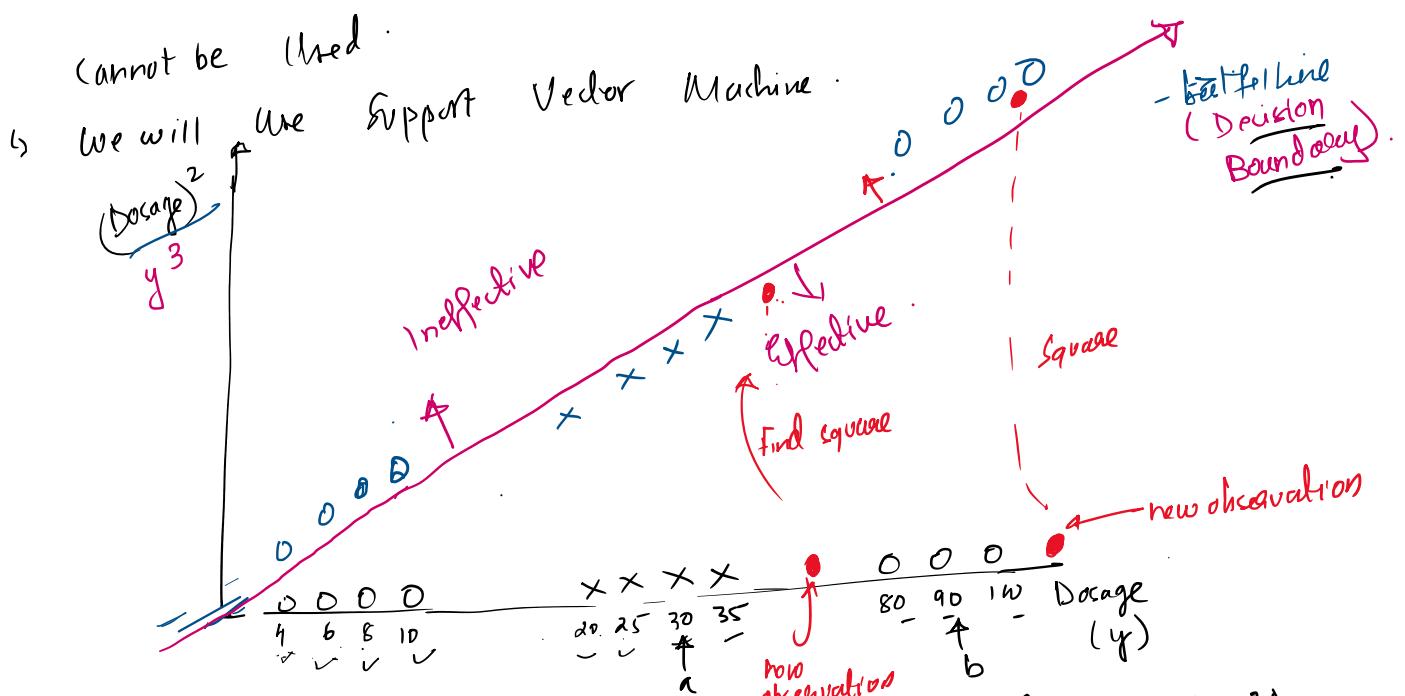
for points outside ϵ Insensitive tube

We have considered up till now Support Vector Classifier
(Soft Margin)
[We were interested in finding Max Margin Classifier]

- * If there are high amount of Overlapping than Soft Margin Classifier will not be efficient
- * In Such Case we will use Support Vector Machine.



→ Here we have lot of Overlapping so Maximum Margin classifier.



→ From above By increasing dimension of data from 1D to 2D and by using Square (as function) we are able to separate the data

* From above we are able to separate the unseparable points and by using Square (as function) we are able to separate the unseparable points.

* We can also use it to predict for new observation
→ For new observation take square of the observation and check if the square lies above (ineffective) or below (effective) the decision boundary.

SVM Main Idea →

- * Start data with low dimension (given).
- * Now move data to higher dimension [Using some function]
- * Find Support Vector Classifier to classify the data into different classes.

Question Arises → Why Square, why not Cube or some other function?

→ How to decide, how to transform.

- * The Function Used be known as Kernel
- * Types of Kernel function in SVM.
 - (1) Polynomial Kernel
 - (2) Radial Basis Kernel.

Note)
if data points in $\underline{1D}$ \rightarrow convex in $\underline{2D} \Rightarrow$ line as Decision boundary
 $\underline{2D} \rightarrow$ convex in $\underline{3D} \Rightarrow$ Plane as Decision Boundary

2)
if Data points in $1D \rightarrow$ Point is Decision Boundary
 $2D \rightarrow$ Line is Decision Boundary
 $3D \rightarrow$ Plane is Decision Boundary

Polynomial Kernel

- * It calculates higher Dimensional Relationship between observations (data points).
 - * The Kernel that transforms 1D to dD is Polynomial Kernel.
 - * It may look like $(a \times b + r)^d$
where a and b are two different observations in dataset
(any two data points).
 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 classifier.
 - * Let a and b be two observations.
 - * Let $r = 1/2$ & $d = 2$ } Note ⇒ value of r and d is determined by cross validation.
- $$\begin{aligned}\Rightarrow (a \times b + r)^d &= (a \times b + 1/2)^2 \\ &= (a \times b + 1/2) \cdot (a \times b + 1/2) \\ &= a^2 b^2 + \frac{1}{2} ab + \frac{1}{2} ab + \frac{1}{4} \\ &= ab + a^2 b^2 + \frac{1}{4} \Rightarrow \text{Polynomial.} \\ &\quad \curvearrowleft \text{can be written as dot product of}\end{aligned}$$

= can be written as dot product of

$$\Rightarrow = \left(\underline{\underline{a}}, \underline{\underline{a}}^2, \frac{1}{2} \right) \cdot \left(\underline{\underline{b}}, \underline{\underline{b}}^2, \frac{1}{2} \right)$$

For data point a

Original value of a
Higher dimension value of a

Original value of b
Value of b in higher dimension

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

Note: $\frac{1}{2}$ is third axis but value same so ignore.

So $(a+b)^d$ is used to get higher dimension "rel" between two data points a & b .

Ex: $\underline{\underline{a}} = \underline{\underline{9}}, \underline{\underline{b}} = \underline{\underline{14}}$

$\underline{\underline{a}} = \underline{\underline{1}} , \underline{\underline{b}} = \underline{\underline{2}}$

\Rightarrow Higher Dimension Relationship betⁿ $a=9, b=14$

$\Rightarrow (\underline{\underline{9}} + \underline{\underline{14}} + \underline{\underline{\frac{1}{2}}})^2 \Rightarrow 126.5^2 \Rightarrow \underline{\underline{16002.25}}$

↑
Relationship representation of
 $a=9, b=14$ in
Higher Dimension

* We can find this Higher Dimension Relationship betⁿ every pair of data points.

VVIMP \rightarrow Kernel Trick

① Kernel function actually never does any transformation in higher dimension.

② Instead it calculates relation betⁿ observations and visualizes them in higher dimension.

visualizes them in higher Dimension.

③ Support Vector Classifier uses this visualization for classification.

④ This is known as Kernel Trick.

$d=1$ \Rightarrow The Polynomial Kernel compute the relationship betⁿ each pair of observation in 1D.

$d=2$ \Rightarrow " " " " " " " " $\in \mathbb{R}^D$.

Best value of d can be found by Cross Validation.

Note Since Polynomial Kernel is $\underbrace{(a \cdot b + 1)^2}_{\text{↑}} = \underbrace{(a, a^2, 1)}_{\text{↑}} \cdot \underbrace{(b, b^2, 1)}_{\text{↑}}$

\rightarrow we actually do not have to Transform to understand higher dimension relationship.

\rightarrow All we need to do is calculate dot product betⁿ each pair of point.

Radial Basis Kernel [VImp]

- * Works in Infinite Dimension.
- * If there are lot of Overlapping in data points classification then Support Vector Classifier cannot be used to linearly separate.
- * One way to deal with overlapping data is to use SVM with Radial Kernel.
- * Radial Kernel uses Radial Basis function!

$$-\gamma \cdot \frac{(a-b)^2}{C}$$

Here a and b are two observations.

$(a-b)^2$ \Rightarrow diff betⁿ the measurement is squared giving us the squared distance betⁿ two observation.

γ \Rightarrow scales the squared distance and thus scales the influence.

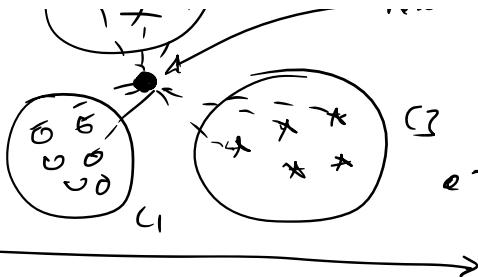
- * Since Radial Basis function finds Support Vector classifier in Infinite Dimension it is not possible to visualise it.

- * When applied to new observation, it behaves like Weighted Nearest Model.

- * Radial Kernel checks the influence of the Nearest



Influence of the Nearest Neighbour on New observation and accordingly makes classification.



- * Let see how radial kernel determines how much influence each observation in training dataset has on classifying new observation.

$$\text{Let } \begin{cases} a = 2.5 \\ b = 4 \end{cases} \Rightarrow \frac{-1}{e} (2.5 - 4)^2 \Rightarrow 0.01 \rightarrow \textcircled{1}$$

$$\gamma = 1$$

$$\begin{aligned} & \begin{cases} a = 2.5 \\ b = 4 \\ \gamma = 2 \end{cases} \Rightarrow \frac{-2}{e} (2.5 - 4)^2 \Rightarrow 0.01 \end{aligned} \quad \left. \begin{array}{l} \text{The } \gamma \text{ scales} \\ \text{the influence.} \end{array} \right\}$$

$$\begin{aligned} & \begin{cases} a = 2.5 \\ b = 16 \\ \gamma = 1 \end{cases} \Rightarrow \frac{-1}{e} (2.5 - 16)^2 \Rightarrow \frac{-1}{e^{13.5}} \approx \text{close to zero} \end{aligned} \quad \hookrightarrow \textcircled{2}$$

From ① & ② If points are close enough then we have high influence.

and if points are far off then we have low influence.

- * Thus to calculate influence b/w two data points, we can plug in the values in Radial Basis fn with appropriate γ .
- OR
- ... in an infinite dimension.

* We get relationship in infinite dimension.

Not from Exam \Rightarrow Let a and b are two observations.

$$\text{Let } \frac{a+b}{2} = \frac{-1}{e^{\frac{a+b}{2}}} (a^2 + b^2 - 2ab) = \frac{-1}{e^{\frac{a+b}{2}}} (a^2 + b^2). e^{\frac{ab}{2}} \quad (1)$$

Taylor Series Expansion \Rightarrow [Allows any f^n to split in infinite sum]

$$f(x) = f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \dots + \frac{f^{(\infty)}(a)}{\infty!}(x-a)^\infty$$

Let $f(x) = e^x$.

$$e^x = e^a + \frac{e^a}{1!}(x-a) + \frac{e^a}{2!}(x-a)^2 + \dots + \frac{e^a}{\infty!}(x-a)^\infty$$

Let $a=0$

$$e^x = 1 + \frac{x}{1!} + \frac{x^2}{2!} + \dots + \frac{x^\infty}{\infty!}$$

$$e^{ab} = 1 + \frac{ab}{1!} + \frac{ab^2}{2!} + \dots + \frac{a^\infty b^\infty}{\infty!}$$

$$e^{ab} = \left(1, \frac{1}{\sqrt{1!}}, \frac{1}{\sqrt{2!}}, \dots, \frac{1}{\sqrt{\infty!}} \right) \cdot \left(1, \frac{1}{\sqrt{1!}}, \frac{1}{\sqrt{2!}}, \dots, \frac{1}{\sqrt{\infty!}} \right)^T \quad (2)$$

from (1) & (2)

$$= \frac{-1}{e^{\frac{a+b}{2}}} \left(1, \frac{a}{\sqrt{1!}}, \frac{a^2}{\sqrt{2!}}, \dots, \frac{a^\infty}{\sqrt{\infty!}} \right) \cdot \left(1, \frac{b}{\sqrt{1!}}, \frac{b^2}{\sqrt{2!}}, \dots, \frac{b^\infty}{\sqrt{\infty!}} \right)$$

$$\text{hel } S = \frac{-1}{e^a} (a^2 + b^2)$$

$$= \left(1, 1, \frac{a}{\sqrt{1}}, \frac{a^2}{\sqrt{2!}}, \dots, \frac{a^\infty}{\sqrt{\infty!}} \right) \cdot \left(1, 1, \frac{b}{\sqrt{1}}, \frac{b^2}{\sqrt{2!}}, \dots, \frac{b^\infty}{\sqrt{\infty!}} \right)$$

We can see Radial Kernel is equal to Dot product that has coordinate for infinite No of Dimension.