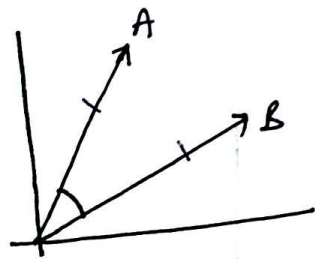# The Similarity Perspective

cosine similarity $= \dfrac{A \cdot B}{\|A\| \, \|B\|}$

$\|A\| = 1 \qquad \|B\| = 2$

$= \boxed{A \cdot B} \rightarrow$ dot product of A and B

\* find the similarity between two vectors we use cosine similarity

$$\max_{\alpha_i} \; \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j \, y_i y_j \, \boxed{(x_i \cdot x_j)} \longrightarrow \text{similarity}$$

$\longrightarrow$ maximize the similarity of sv based on main sign

# Kernel SVM

$$\max_{\alpha_i} \; \sum_{i=1}^{n} \alpha_i - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_i \alpha_j \, y_i y_j \, k(x_i, x_j) \longrightarrow \boxed{\text{Kernel SUM}}$$

$k(x_i, x_j) \rightarrow$ kernel $\rightarrow$ similar between $x_i$ and $x_j$

$\longrightarrow \boxed{x_i \cdot x_j} \longrightarrow \boxed{x_i \cdot x_j}$

$\longrightarrow$ Linear SVM

$\boxed{\text{Polynomail}} \qquad \boxed{\text{RBF}} \qquad$ Creative Version

# Polynomial Kernel



$$x_1^2 + x_2^2 \longrightarrow$$

|  | $x_1$ | $x_2$ |
|---|---|---|
| $x_i \rightarrow$ | $x_{11}$ | $x_{12}$ |
| $x_j \rightarrow$ | $x_{21}$ | $x_{22}$ |

$$\boxed{x_{11}\, x_{21} + x_{12}\, x_{22}}$$

$$k(x_i, x_j) = (\delta + x_i \cdot x_j)^d \qquad \delta = 1 \qquad d = 2, 3$$

$$= \boxed{1 + x_{11}^2\, x_{21}^2 + x_{12}^2\, x_{22}^2 + 2x_{11}\, x_{21} + 2x_{12}x_{21} + 2x_{11}\, x_{21}\, x_{12}x_{22}}$$

$\downarrow$ Polynomial term

## Kernal Trick

$x_1$ and $x_2$    not    convert    into    other form ( 3d or 4d... )

$x_1$ and $x_2$    use    Polynomial term    in

current form.

## The trick

$$1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11} x_{21} + 2x_{12}x_{22} + 2x_{11} x_{21} x_{12} x_{22}$$

dot product of 2 Vector

$$\left[\begin{array}{cccccc} 1 & x_{11}^2 & x_{12}^2 & \sqrt{2}\,x_{11} & \sqrt{2}\,x_{12} & \sqrt{2}\,x_{11}\,x_{21} \end{array}\right] \xrightarrow{\ x_i\ } \text{6d Vector}$$

$$\left[\begin{array}{cccccc} 1 & x_{22}^2 & x_{21}^2 & \sqrt{2}\,x_{21} & \sqrt{2}\,x_{22} & \sqrt{2}\,x_{12}\,x_{22} \end{array}\right] \to \text{6d vector} \to x_j$$

1 method $\to$

$$\begin{bmatrix} x_i\,(x_{11}\ x_{12}) \xrightarrow{ft} x_i\,(6d) \\ x_j\,(x_{21}\ x_{22}) \to x_j\,(6d) \end{bmatrix} \Large{>} \normalsize x_i' \circ x_j' \to \text{expression}$$

2 method $\quad \dfrac{x_i}{x_j} \Large{>} \normalsize k(x_i, x_j) \to$ expression

✗ 1 method occupy much space to store the vector

✔ 2 method not change into higher dimension just use polynomial term calculate expression.

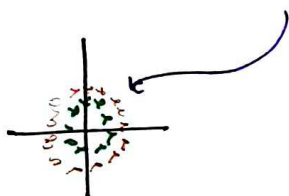Let take 10 dimensional data if we solve with method 1 then occupy large space bcz change into higher dimension.

but in 2 method not change into higher dimension and not occupy any space.

# What about the other Polynomial terms

$$1 + X_{11}^2 X_{21}^2 + X_{12}^2 X_{22}^2 + 2X_{11} X_{21} + 2X_{12} X_{22} + 2X_{11} X_{21} X_{12} X_{22}$$

Circular shape

Other shape ( Conic section)

# RBF Kernel

Radial Basic function { Normal distribution }

→ Popular

→ Best out of the box kernel

→ Powerful

euclidian distance

$$k(X_i, X_j) = \left\{ e^{\dfrac{-\|x_i - x_j\|^2}{2\sigma^2}} \right\}$$

$k \propto \dfrac{1}{distance}$

similarity

distance

$$c^{-\gamma \| x_i - x_j \|^2} \qquad \boxed{\gamma = \dfrac{1}{2}\sigma^2}$$

gamma

→ hyperparameter

$$k(x_i, x_j) = e^{-\frac{dist^2}{2e^2}}$$

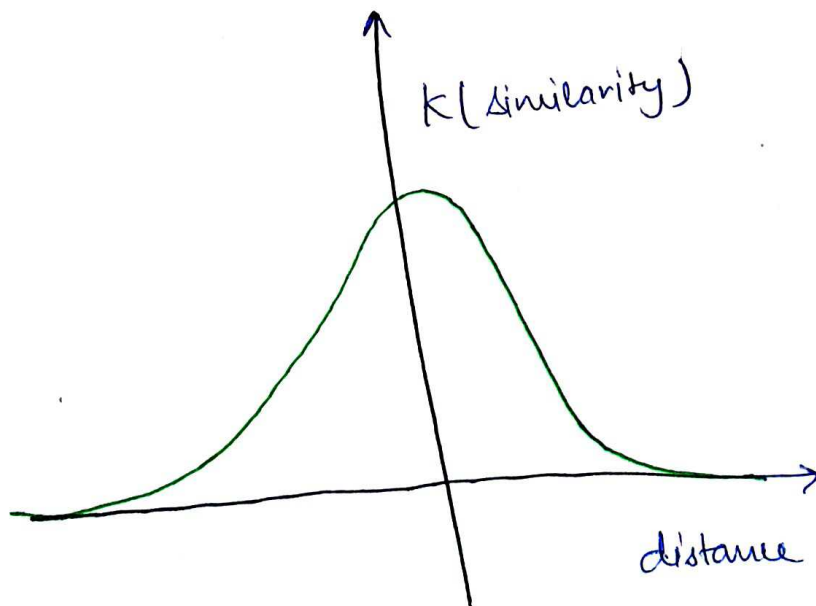$$k \propto \frac{1}{distance}$$

↑

Situation

- **Non- linear Transformations** : The RBF kernel enables the use of non-linear transformations, which can map the original feature space to a highest-dimensional space when the data becomes linearly seperable. This is particularly useful for problems where the decision boundary is not linear.

- **Local Descision**: Unlike some other kernels, the RBF kernel make "local" delisions. That is, the effect of each data point is limited to a certain region around that point. This can make the model more robust to outliers and create complex descision boundary.

- **Flexiblity** : The RBF kernel has a parameter γ (related to the standard deviation of the Gaussian distribution) that determines the complenity of the decision boundary. By tuning this parameter, we can adjust the trade-off between bias and variance, allowing for a flexible range of decision boundaries.

- **Universal Approximation Property:** The RBF Kernel has a property known as the universal approximate property, meaning it can approximately any continuous function to a certain degree of accuracy given enough data points. This makes it highly versatile and capable of modeling a wide variety of relationship in data.

- **General purpose:** The RBF kernel does not makes any strong assumption about the versatile, general-purpose Kernel.
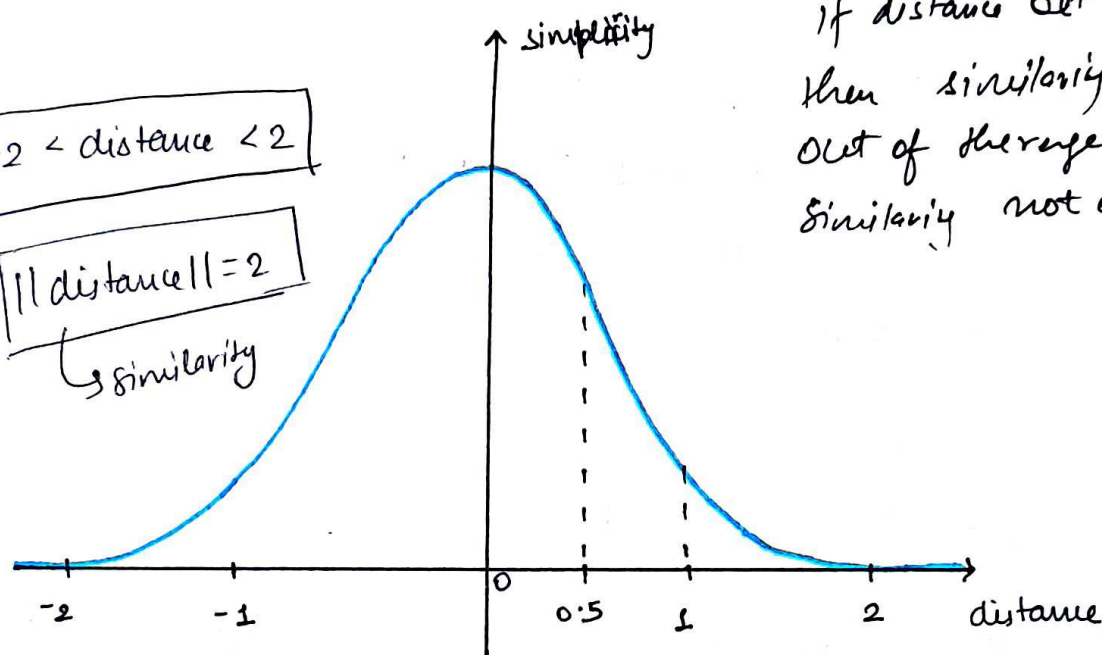
## Local Decision Boundary

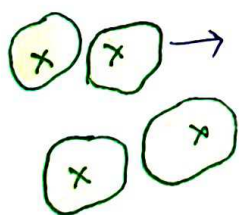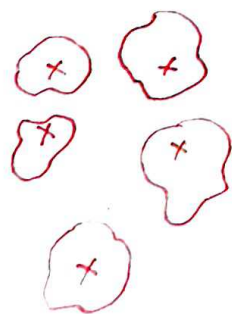$$\left[ e - distance^2 \right] \rightarrow X$$

$-2 <$ distance $< 2$

$\|$ distance $\| = 2$
$\quad \hookrightarrow$ similarity

if distance betⁿ -2 and 2
then similarity exist
out of the range (-2 and 2)
similarity not exist.

similarity ↑

-2    -1    0    0·5    1    2    distance

$e^{-\|x_i - x_j\|^2}$ → similarity and distance is
exponential decrease.

-2 < distance < 2

if any point in the
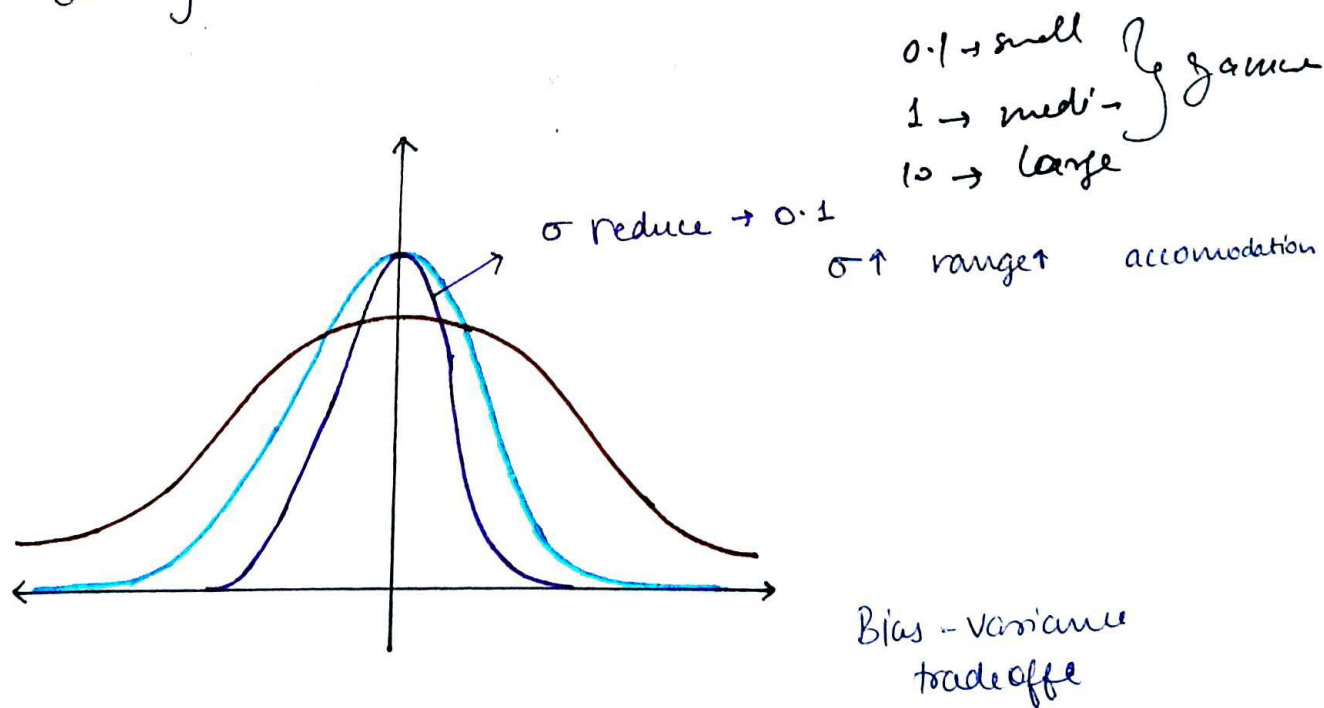area then similar
otherwise not
similar.

# Effect of Gamma

The parameter $\gamma$ in the Radical Basis function (RBF) kernel of a support vector Machine (SVM) is a hyperparameter that determines the spread of the kernel and therefore the decision region.

The effect of $\gamma$ can be summarized as follows:

- If $\gamma$ is too large, the exponential will decay very quickly, which means that each datapoint will only have an influence in its immediate vicinity. The result in a more complex decision boundary, which might overfit the training data.

- If $\gamma$ is too small, the exponential will decay slowly, which means that each data point will only have an ~~influence~~ in ~~its~~ ~~immediate~~ ~~vicinity~~ wide range of influence. The decision boundary will therefore be smoother and more simplistic, which underfit the training data.

In a sense, γ in the RBF kernel plays a role similar to that of the inverse of the regularization parameter. It controls the trade-off between bias (underfitting) and variance (~~offer~~ overfitting). High γ values can lead to high variance (overfitting) due to more flexibility in shaping the decision boundary, while low γ values can lead to high bias (underfitting) due to a more rigid, simplistic decision boundary.

Tuning the γ parameter using cross-validation or a similar technique a crucial step when training SVMs with an RBF kernel.

$0.1 \to$ small $\left.\begin{array}{l} \\ \end{array}\right\}$ γ value
$1 \to$ medi-
$10 \to$ large



σ reduce $\to 0.1$

$\sigma \uparrow$ range $\uparrow$ accommodation

Bias - variance trade off

$\sigma \downarrow$ or $\gamma \uparrow \to$ overfitting

$\sigma \to$ hyperparameter

$\sigma \uparrow$ or $\gamma \downarrow \to$ underfitting

$\gamma \downarrow \to$ locality $\uparrow$   $\sigma \propto \dfrac{1}{\gamma}$

$\gamma \uparrow \to$ locality $\downarrow$

# Relation Between RBF and Polynomial kernel

Infinite Dimensional Mapping: The RBF kernel implicity maps input data to an infinite dimensional feature space, which allows for even greater flexibility in forming decision boundaries.
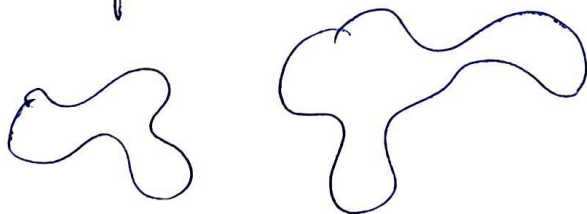
In polynomial

$$X_1 \quad X_2 \xrightarrow{\text{degree} \to 2} X_1 \quad X_2 \quad X_1^2 \quad X_2^2 \quad X_1 X_2$$

In RBF $\to$ $\infty$ dim vector

$X_1 \quad Y_2$ $\Bigg\downarrow$ $\infty$ dim vect

$$X_1 \quad X_2 \quad X_1^2 \quad X_2^2 \quad X_1 X_2 \to d=2$$

$$X_1^3 \quad X_1^2 \quad X_1^2 X_2 \quad X_2^2 X_1 \to d=3$$

$$X_1^4 \quad X_2^4 \quad X_1^3 X_2^1 \quad X_1^2 X_2 \quad X_2^3 X_1^1 \to d=4$$

$$\vdots$$

$$d = \infty$$

\* RBF $\to$ Because of $\infty$ dimension vector RBF make any type of decision boundary.

$$K(X_i, X_j) = e^{-\|X_i - X_j\|^2 \over 2\sigma^2} \qquad \text{let} \quad \sigma = 1 \quad \sigma^2 = 1$$

$$= e^{-\|X_i - X_j\|^2 \over 2}$$

$$= e^{-{(x_i - x_j)^T (X_i - X_j) \over 2}} \qquad = e^{-{(x_i^T - x_j^T)(X_i - X_j) \over 2}}$$

$$X_i = \begin{bmatrix} X_{11} & X_{12} \end{bmatrix}$$

$$X_j = \begin{bmatrix} X_{21} & X_{22} \end{bmatrix}$$

$$= e^{-{[x_i^T x_j - x_j^T x_j - x_j^T x_i + x_j^T x_j] \over 2}}$$

$$= e^{-{[x_i^T x_i + x_j^T x_j - x_i^T x_j - x_j^T x_i] \over 2}}$$

$$X_i^T X_j =$$

$$X_j^T X_i =$$

$$= \underbrace{e^{-{1 \over 2}[x_i^T x_i + x_j^T x_j]}}_{\text{constant}} \; e^{X_i^T X_j}$$

$$= C e^{1 + X_i^T X_j - 1}$$

$$= C \, e^{1 + X_i^T X_j} \; e^{-1}$$

$$= C' \sum_{k=0}^{\infty} {(1 + X_i^T X_j)^k \over k!} \qquad \left[\because e^x = \sum_{k=0}^{\infty} {x^k \over k!}\right]$$

$$\hookrightarrow \; C' \sum_{k=0}^{\infty} k_{poly} {(X_i, X_j)^k \over k!}$$

$$C'\left[1 + {(1 + X_i^T X_j) \over 1!} + {(1 + X_i^T X_j) \over 2!} + \dots {(1 + X_i^T X_j)^{\infty} \over \infty!}\right]$$

proved

$$RBF = \sum_{i=0}^{\infty} k_{poly}(x_i, x_j)$$

# Custom Kernels

1. **String kernel:** These are used for classifying text or sequences, where the input data is not numerical. String kernels measure the similarity between two strings. For example, a simple string kernel might count the number of common substring between two strings.

2. **Chi-Square kernel:** This kernel is often used in computer vision problems, especially for histogram comparison. It's defined as $k(x,y) = \exp(-\gamma x^2(x,y))$, where $x^2(x,y)$ is the chi-square distance between the histogram x and y.

3. **Intersection kernel:** This is another kernel commonly used in computer vision, which compute the intersection between two histograms (or generally non-negative feature vectors).
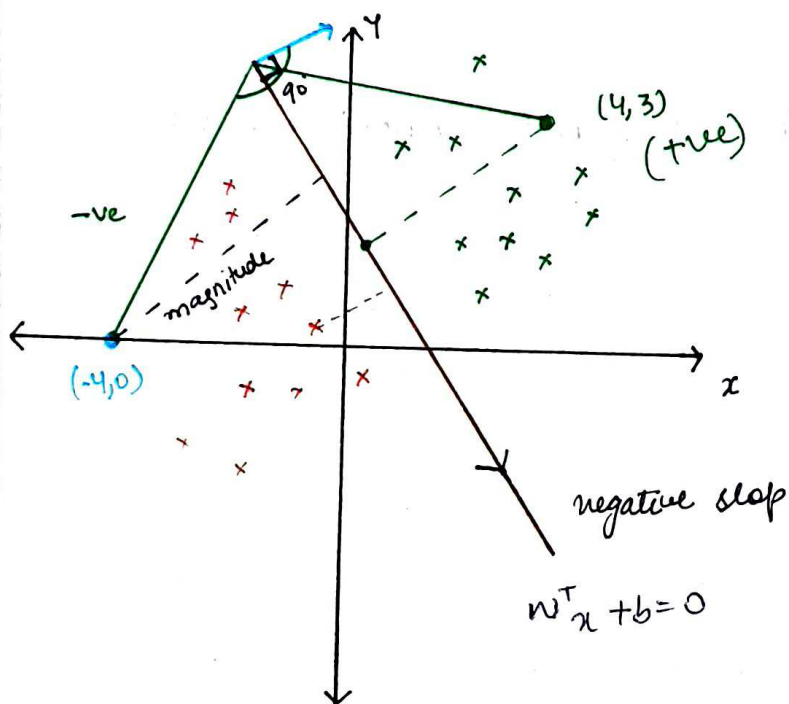
4. **Hellinger's kernel:** Hellinger's kernel, or Bhattacharya kernel, is used for company probability distribution and is popular in image recognition tasks.

5. **Radial basis function network (RBFN) kernels:**

These are similar to the standard RBF kernel, but the centers and width of the RBFs are learned from the data, rather than being fixed priori.

6. **Spectral kernels:** These kernels use spectral analysis techniques to compare data points They can be particularly useful for dealing with cyclic or periodic data.

$$y = mx + c$$

$$h\theta(x) = \theta_0 + \theta_1 x_1$$

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \theta_3 x_3$$
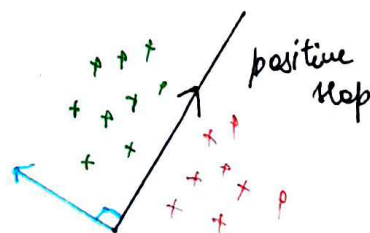
$$y = b + [w_1 x_1 + w_2 x_2 + w_3 x_3]$$

$$W = \begin{bmatrix} w_1 \\ w_2 \\ w_3 \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

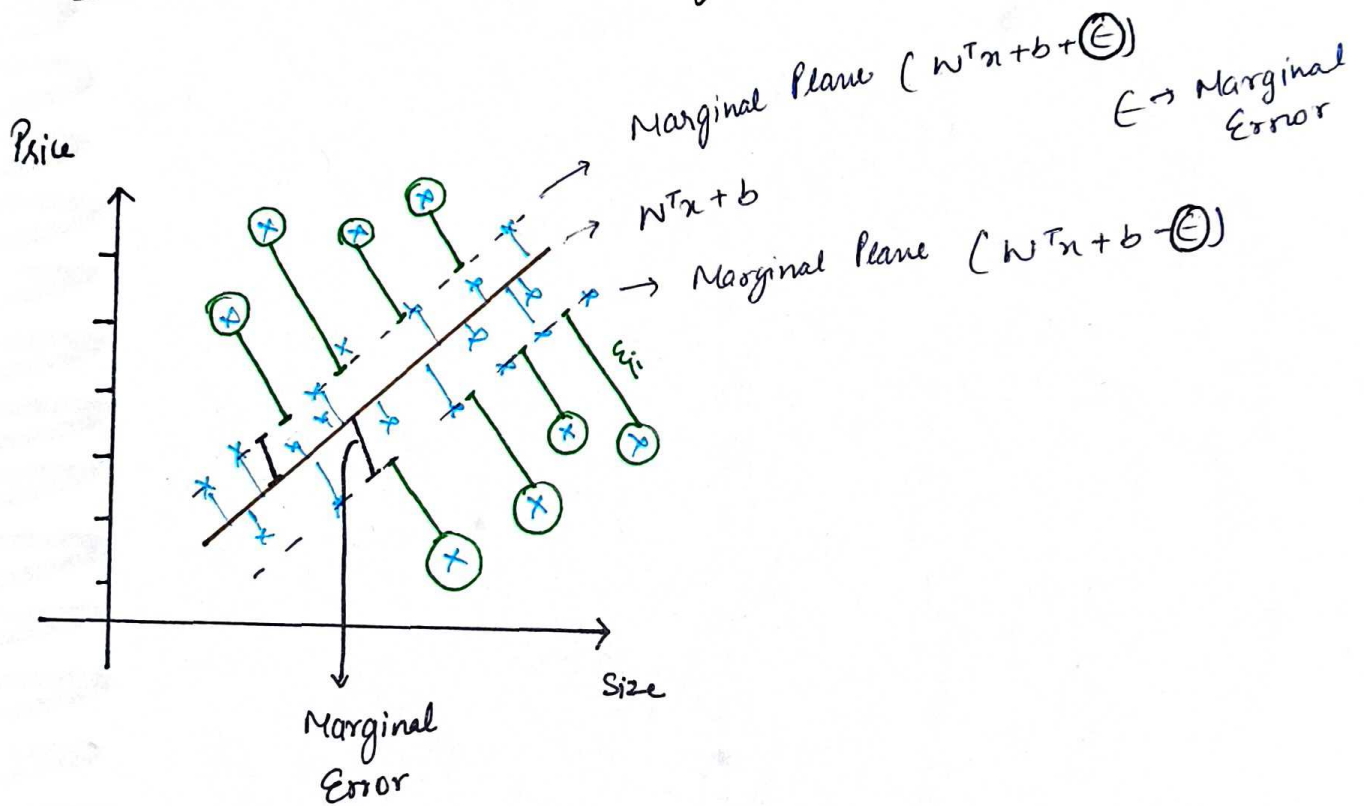**#** If angle of point is more than $90°$ is ~~positive~~ negative point

$$W^T = \begin{bmatrix} w_1 & w_2 & w_3 \end{bmatrix} \qquad x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

$$W^T x = 0$$

**#** If angle of point is less than or equal to $90°$ is positive point.

# SVR ( Support Vector Regressor )



Price

Marginal Plane $(w^Tx + b + \epsilon)$

$\epsilon \rightarrow$ Marginal Error

$w^Tx + b$

Marginal Plane $(w^Tx + b - \epsilon)$

$\xi_i$

Size

Marginal Error

## Cost function

$$\underset{w, b}{\text{Min}} \quad \frac{||W||}{2} + \boxed{C \sum_{i=1}^{m} \xi_i} \Rightarrow \text{Hinge}$$

## Constrain:

$$\text{Error} \in |y_i - w^Tx_i| \leq \epsilon$$

\* points under the margin or distance between points and best fit line and error always less than margin

$$\text{Error} \in |y_i - w^Tx_i| \leq \epsilon + \xi_i$$

\# some points outside the margin do, we use $\xi_i$ and add distance between the margin and outside point.