

Components of a neural network

Activation Function

Unit Step Activation Function

* Linear



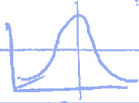
* Saturated Linear



* Hyperbolic Tangent

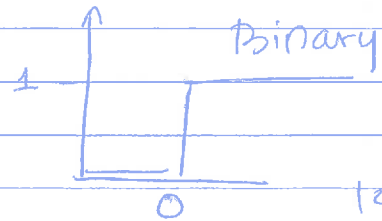


* Gaussian



↓

Radial Basis Function (RBF)

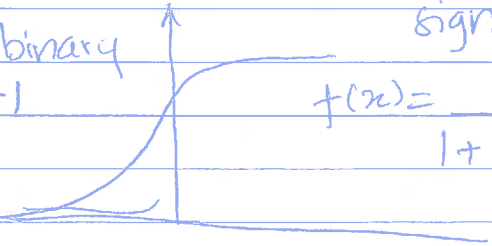


logistic sigmoid.

$$f(x) = \frac{1}{1 + e^{-x}}$$

Not binary

any 0-1



Network Topography

Single-layer network → linearly separable data

- Feed forward: → Multilayer Perceptron (MLP)
- Deep Neural Network (DNN)
- Recurrent Network (both dirⁿ flow of signal)
Delay → understand seq of data over time.



Use fewest nodes that result in adequate performance

STRENGTH

- * Classification / numeric
- * modelling complex patterns better than many other algos
- * makes few assumptions about the data.

WEAKNESS

- * Computationally intensive
- * Prone to overfitting.
- * Complex Black Box Model that are difficult to understand.

Learning Rate : Rate of gradient descent

↑ rate ↓ training time.

Forward then Backward to check error - Epoch

Neural Networks

③

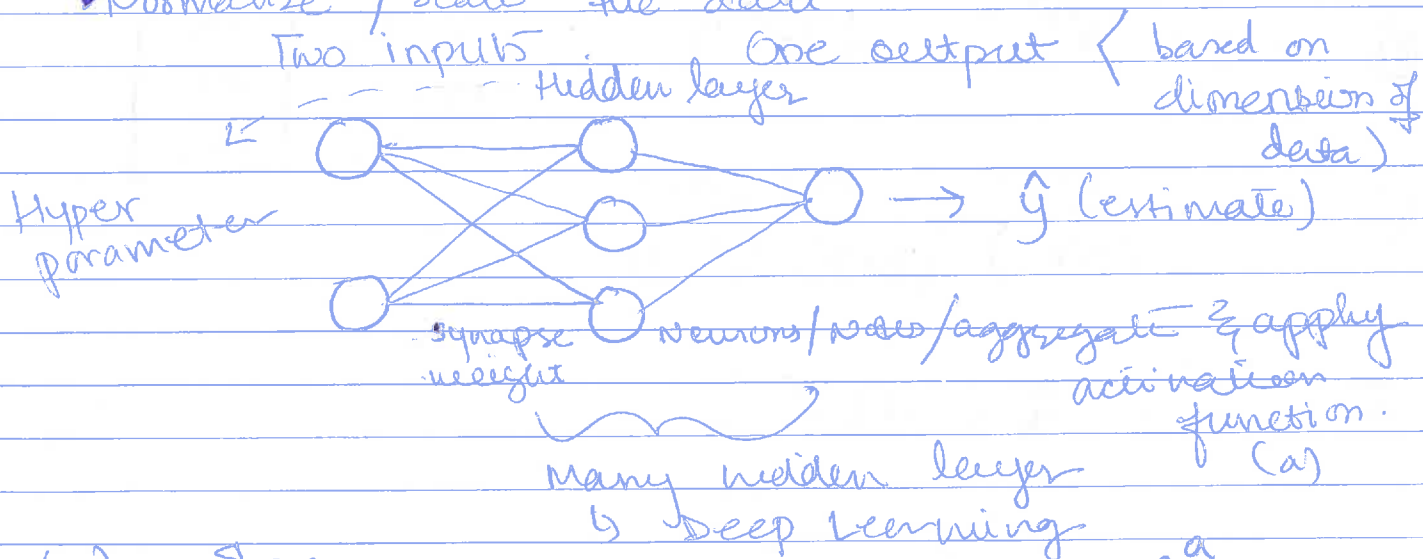
	hours study X hours sleep	test (y) score	
Train	(3, 5)	75	supervised regression if letter grade then classification
	(5, 1)	82	
	(10, 2)	93	
Test	(8, 3)	?	

Usually good with Image recognition, music, text classification etc.

NPL

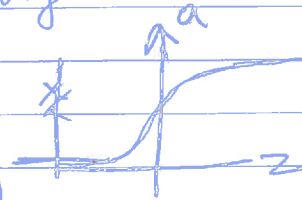
* Need to account for differences in scale / units of the data.

↳ Normalize / Scale the data



$$(z) = \sum x_i$$

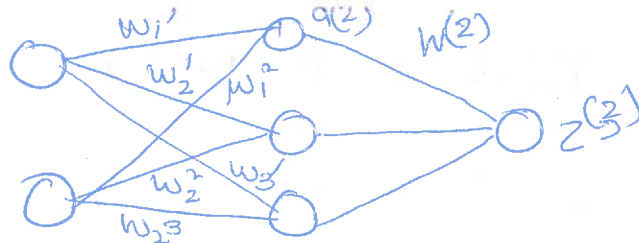
(a) activation function = $\frac{1}{1 + e^{-z}}$ (Sigmoid fn)



Structure & behaviour of network is defined by Hyperparameter not updated as we train the network, what is learned are the weights



$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix}$$



④

$$\begin{bmatrix} 3 & 5 \\ 5 & 1 \\ 10 & 2 \end{bmatrix} \begin{bmatrix} w_{11}^{(1)} & w_{12}^{(1)} & w_{13}^{(1)} \\ w_{21}^{(1)} & w_{22}^{(1)} & w_{23}^{(1)} \end{bmatrix} = \begin{bmatrix} 3w_{11} + 5w_{21} & 3w_{12} + 5w_{22} \\ 5w_{11} + 1w_{21} & 5w_{12} + 1w_{22} \\ 10w_{11} + 2w_{21} & 10w_{12} + 2w_{22} \end{bmatrix}$$

① --- X $W^{(1)}$ $Z^{(2)}$
 (3×2) (2×3) (activity of 2nd layer)
 (3×3)

Next we apply the activation function f on each value in $Z^{(2)}$

$$a_{11} = \frac{1}{1 + e^{-z_{11}}}$$

Sigmoid function in r
 (3×3)

② --- $a^{(2)} = f(z^{(2)})$

$$W^{(2)} = \begin{bmatrix} w_{11} \\ w_{21} \\ w_{31} \end{bmatrix}$$

3x1

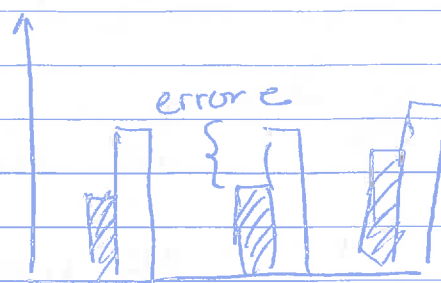
③ --- $z^{(3)} = a^{(2)} W^{(2)}$
 (3×1) (3×3) (3×1)

④ --- $\hat{y} = f(z^{(3)})$

< Next we will train >

GRADIENT DESCENT

What we predicted vs Actual



$$\text{error} = y - \hat{y}$$

$$\begin{aligned} \text{cost function } J &= \sum \frac{1}{2} (e_1^2 + e_2^2 + e_3^2) \\ &= \sum \frac{1}{2} e_i^2 \end{aligned}$$

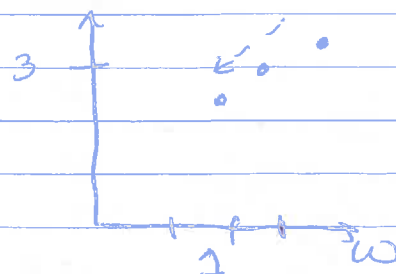
Training \rightarrow minimizing the cost function (overall error)

cost function \rightarrow ^{inputs} examples & weights
 we cannot change the inputs
 so we change the weight to
 minimize cost

"CURSE OF DIMENSIONALITY"

to avoid calculating weights (all) we look at w and calculate cost

For example in 1D, if
 if know which way the
 error is decreasing we can
 go in that direction
"Numerical Estimation"

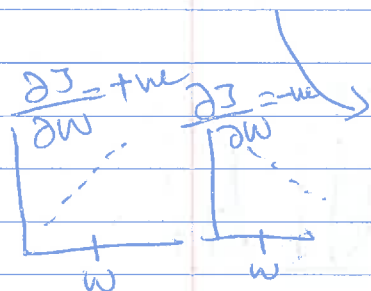


BUT...

⑥

Let's look at our equations

$$J = \sum \frac{1}{2} (y - f(f(xw^{(1)}) * w^{(2)}))^2$$



Rate of change of J w.r.t $w^{(1)}$ & $w^{(2)}$ will tell us which way is down

We consider partial derivative \rightarrow one weight at a time

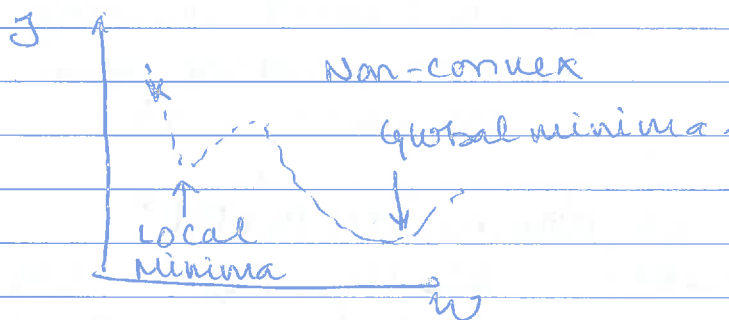
$$\frac{\partial J}{\partial w}$$

Take steps iteratively till cost stops decreasing

\rightarrow GRADIENT DESCENT

Works well when high dimension data.

Problem!



By making the cost square we

make this convex



Stochastic
STOCHASTIC
GRADIENT DESCENT

one example at a time to calculate $\frac{\partial J}{\partial w}$ and move in that direction

Sometimes one at a time does not get us stuck in local minima

7

Batch Gradient Descent (adds the cost of one at a time)

$$\frac{\partial J}{\partial W_2} = \frac{\partial \sum \frac{1}{2} (y - \hat{y})^2}{\partial W_2}$$

$$= \cancel{\sum} \frac{\partial \frac{1}{2} (y - \hat{y})^2}{\partial W_2}$$

$$= \cancel{\sum} 2 \times \frac{1}{2} \cancel{(y - \hat{y})} (y - \hat{y}) \frac{\partial \hat{y}}{\partial W_2} \leftarrow \hat{y} = f(z^{(3)})$$

$$= \cancel{\sum} \underset{(\text{drop})}{-} (y - \hat{y}) \frac{\partial \hat{y}}{\partial z^3} * \frac{\partial z^3}{\partial W_2}$$

$$\begin{cases} f(z) = \frac{1}{1 + e^{-z}} \\ f'(z) = \frac{e^{-z}}{(1 + e^{-z})^2} \end{cases}$$

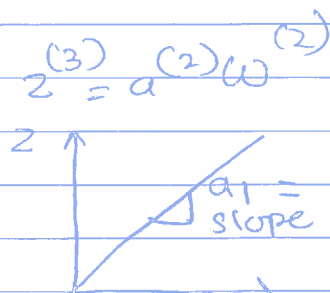
$$\frac{\partial J}{\partial W_2} = \sum \underset{(\text{drop})}{-} (y - \hat{y}) f'(z^3) \frac{\partial z^3}{\partial W_2}$$

$$\begin{matrix} \downarrow & \begin{bmatrix} 3 \times 1 \end{bmatrix} & \begin{bmatrix} 3 \times 1 \end{bmatrix} & \begin{bmatrix} a^2 \\ 3 \times 3 \end{bmatrix} \\ & & g^{(3)} & \end{matrix}$$

when we have to do the summation we can transpose a^2

$$a^T g^{(3)} = \frac{\partial J}{\partial W_2}$$

$$g^{(3)} = -(y - \hat{y}) f'(z^3)$$



if we look W_{11} at one synapse $z^3 W$ have a linear relation

8

$$\frac{\partial J}{\partial w_1} = \sum_1 \frac{1}{2} (y - \hat{y})^2$$

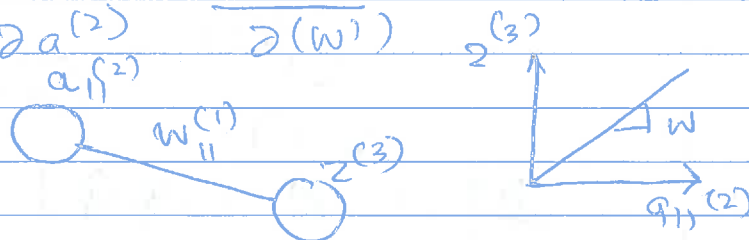
$$= -(y - \hat{y}) \frac{\partial \hat{y}}{\partial w^{(1)}}$$

$$= -(y - \hat{y}) \frac{\partial \hat{y}}{\partial z^3} \cdot \frac{\partial z^3}{\partial w^{(1)}}$$

$$= -(y - \hat{y}) \cdot f'(z^{(3)}) \frac{\partial z^3}{\partial w^1}$$

$$= \delta^{(3)} \frac{\partial z^{(3)}}{\partial (w^1)}$$

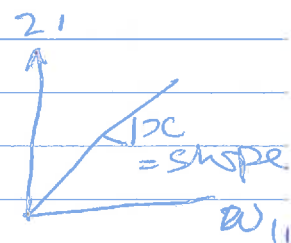
$$= \delta^{(3)} \frac{\partial z^{(3)}}{\partial a^{(2)}} \cdot \frac{\partial a^{(2)}}{\partial (w^1)}$$



$$= \delta^{(3)} (w^{(2)})^T \frac{\partial a^{(2)}}{\partial w^{(1)}}$$

$$= \delta^{(3)} (w^{(2)})^T \frac{\partial a^{(2)}}{\partial z^{(2)}} \cdot \frac{\partial z^{(2)}}{\partial w^{(1)}}$$

$$= \delta^{(3)} (w^{(2)})^T f'(z^{(2)}) \cdot \frac{\partial z^{(2)}}{\partial w^{(1)}}$$



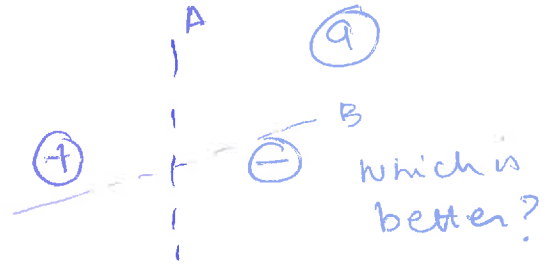
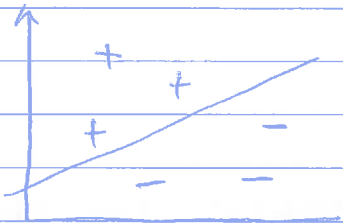
$$\frac{\partial J}{\partial w_1} = \underbrace{x^{(T)} \delta^{(3)} (w^{(2)})^T f'(z^{(2)})}_{\delta^{(2)}}$$

Support Vector Machines.

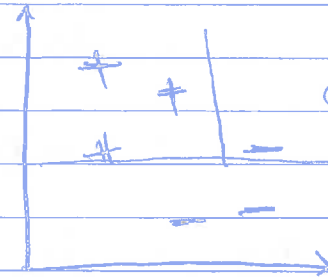
(Supervised)

Draw boundaries.

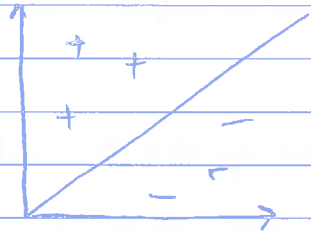
neural net



trees.
axis parallel splits

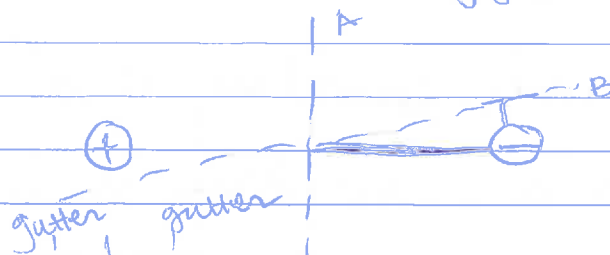


SVM



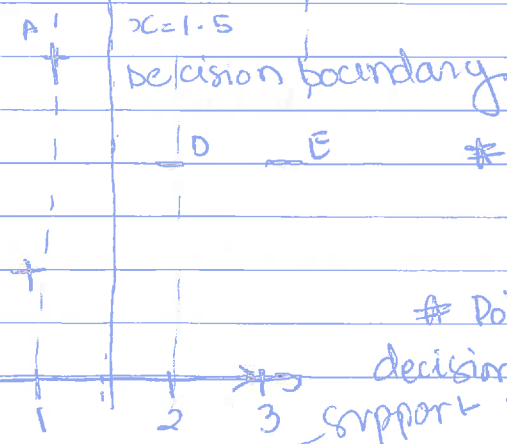
SVM separates the data with the largest possible margins. (gutter)

- SVM is a numerical classifier \rightarrow data has numerical features.
- Draws a single decision boundary to maximize the margin width or gutter. (b/w + & - points)
- It is a binary classifier. (+/-)
w/ a neural net classify as (1/0)



Distance to training point is max hence

A, B, D
ne support vectors because
if we move any
the boundary would be
different



* Cannot be in the gutter but can be on the gutter.

* Points necessary to define the decision boundary are called SV support vectors

4 training points

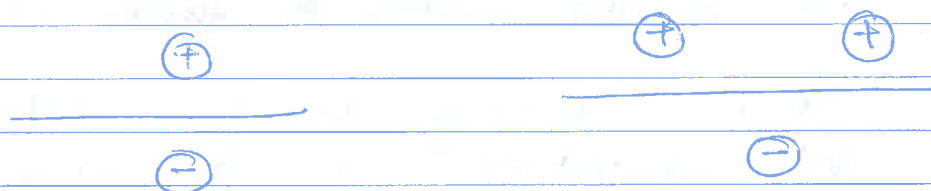
2 classes of points

3 support vectors

Margin is the area between the two gutters.
Margin width is the distance from one gutter to another.
Margin width = 1

We can add points anywhere else \rightarrow no change

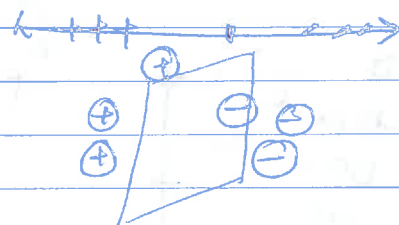
add points on the gutter \rightarrow Are those support vectors
so if we add the points such that if we
remove them then the boundary would
change \rightarrow the support vectors.



in 2D our boundary is a line (2/3 sv)

in 3D our boundary is a point (2 DV)

in 3D \rightarrow plane

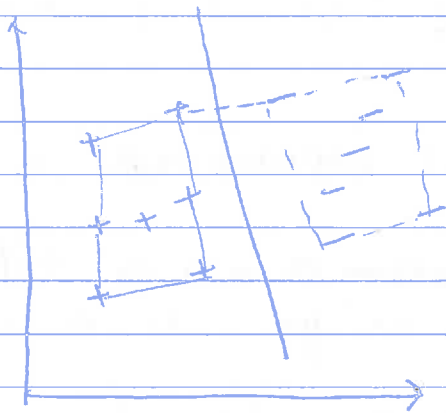


Always need at least 2 SVs.

Where does the decision boundary go?

We use the method called CONVEX HULL then we search for Maximum Margin Hyperplane (MMH). MMH selected is the one that creates the greatest separation.

* Identifying SVs allows to store the classification model very compactly even in case of large # of features.



Outer boundary is known as Convex Hull
 - Find the shortest line between the two convex hulls.
 - Bisector of that line is the Decision Boundary.

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

weights $\vec{w} = \{ \vec{w}_1, \vec{w}_2, \dots, \vec{w}_n \}$

$b = \text{bias} \{ \text{intercept term} \}$

goal is to find weights, such that

$$\left. \begin{array}{l} \vec{w} \cdot \vec{x} + b \geq +1 \\ \vec{w} \cdot \vec{x} + b \leq -1 \end{array} \right\} \text{weights that specify two hyperplanes.}$$

Distance b/w two points = $\frac{2}{\|\vec{w}\|}$ \rightarrow euclidean norm (distance from origin to vector)

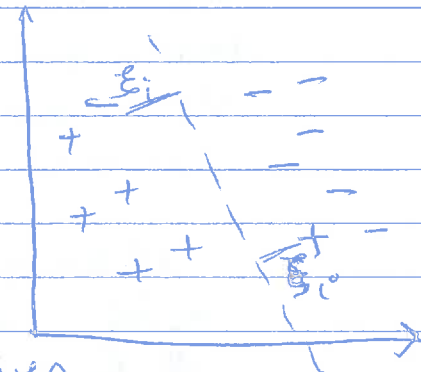
to max distance, we need to min $\|\vec{w}\|$

$$\min \frac{1}{2} \|\vec{w}\|^2$$

s.t. $y_i (\vec{w} \cdot \vec{x}_i - b) \geq 1$, $\forall \vec{x}_i \neq$ subject to each y_i is correctly classified.
 subject to $y_i \in \{+1, -1\}$ for all

What if the data is not linearly separable?

⇒ Use of slack variable
points are allowed to be on the wrong side by defining a soft margin using slack

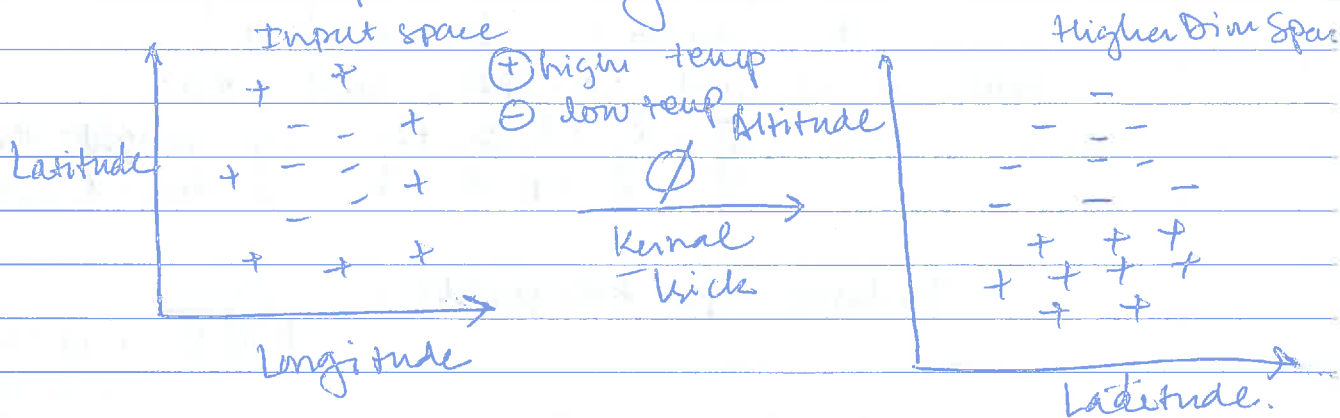


In this case a cost is applied to data points that violate the constraint ξ_i the algorithm tries to minimize the COST instead of maximizing the MARGIN

Adjustment of the "COST PARAMETER" is hence important

↑ CP ⇒ algo ~~tries~~ strives for 100% separation
↓ CP ⇒ algo strives for wider overall margin

⇒ Use of Kernel trick <non-linear kernels>
Ability to map the problem into a higher dimension space using kernel trick



trick is to construct new features that express mathematically relationship b/w measured characteristics

↳ SVM learns concepts not explicitly measured in the original data.

Strength

- ① Classification or numeric prediction.
- ② Not prone to noisy data.
- ③ Not prone to overfitting.
- ④ May be easier to use than neural networks
- ⑤ High accuracy & wins in Data Mining competitions

Weakness

- ① Finding best model requires testing of various combos of Kernel & Model parameters
- ② Can be slow to train
- ③ Results in a complex black box that is difficult to understand.

KERNEL FUNCTIONS

$$K(\vec{x}_i, \vec{x}_j) = \phi(\vec{x}_i) \cdot \phi(\vec{x}_j)$$

ϕ - mapping data in another space.
 Dot product
 return single scalar number.

* LINEAR KERNEL

$$K(\vec{x}_i, \vec{x}_j) = \vec{x}_i \cdot \vec{x}_j$$

* POLYNOMIAL KERNEL (simple non linear transformation)

$$K(\vec{x}_i, \vec{x}_j) = (\vec{x}_i \cdot \vec{x}_j + 1)^d$$

* SIGMOID KERNEL (similar to neural network activation function)

$$K(\vec{x}_i, \vec{x}_j) = \tanh(\underbrace{\kappa \vec{x}_i \cdot \vec{x}_j}_{\text{kappa}} - \underbrace{\delta}_{\text{delta}})$$

* GAUSSIAN RBF Kernel

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

Example next ...