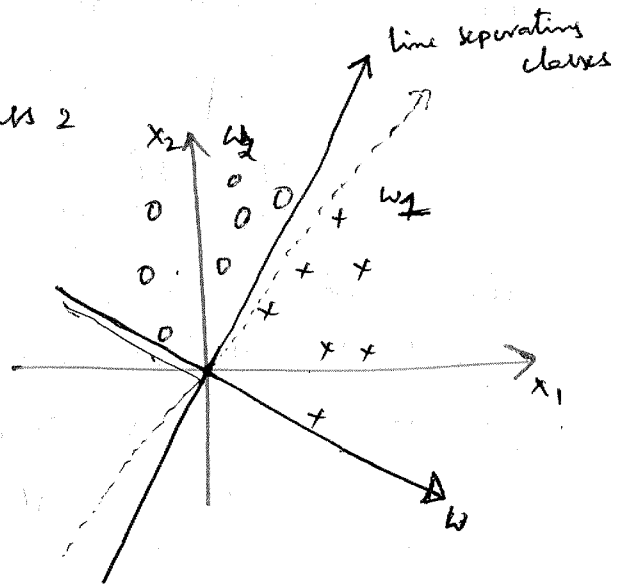


Perceptron learning algorithm:

(1)

Consider Samples of class 1 & class 2

Let there samples are linearly separable. The idea here is to search for the hyperplane that linearly separates two class data samples such that



(if) $w^T x > 0$ for class 1

and $w^T x < 0$ for class 2.

④ one such direction 'w' for the given set of samples of class '1' & '2' is shown. (a) Note that direction of 'w' will be \perp to line separating two classes. (b) Also note that this line separating two classes is not unique. (i.e. even the dotted line also separates class (1) & (2)).

* Let the misclassified ~~with~~ samples of class '1' & class '2' constitute the total cost in misclassification, which is

given by
$$J(w) = \sum_{x \in \text{misclassified } w_1 \& w_2} \delta_x w^T x$$

$$\delta_x = -1 \text{ for class 1}$$

$$1 \text{ for class 2}$$

Note:

This can be extended to lines not passing through origin i.e.

$w^T x + w_0 = 0$ by using $w' = [w \ w_0]$; $x' = [x, 1]$

and comparing $w'^T x' > 0$ for class 1
 $w'^T x' < 0$ for class 2.

The cost function associated with linear discriminant function 'w' is given by

$$J(w) = \sum_{x \in \text{misclassified}} \delta_x w^T x \quad J(w) = \sum_{x \in \text{misclassified}} \delta_x w^T x \rightarrow (1)$$

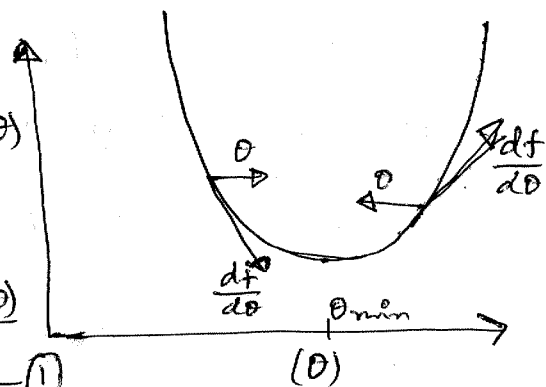
~~But~~ note that $J(w)$ is a piecewise linear function but has discontinuous gradients when set of misclassified vectors 'x' change. Thus the minima of $J(w)$ can not be obtained by direct differentiation.

*) But nevertheless the gradient descent approach can be applied to find minimum of $J(w)$.

⇒ Gradient descent can be applied to ~~find~~ ^{reach} minima ~~by~~ iteratively.

for example we update 'θ' value $f(\theta)$ with following scheme

$$\theta(\text{new}) = \theta(\text{old}) - p \frac{df(\theta)}{d\theta} \quad (1)$$



⇒ i.e. we move in the direction ~~of~~ opposite to (sign) of gradient.

It can be shown that for 'p' satisfying some properties (i.e. $0 < p < 1$ & is decreasing real no. from iteration to iteration)

eq. (1) converges to θ_{\min} .

Now by applying gradient descent to eq. (1) we get

$$w(t+1) = w(t) - p(t) \frac{dJ(w)}{dw}$$

Since $\frac{\partial \delta_x w^T x}{\partial w} = x$

or

$$w(t+1) = w(t) - p(t) \cdot \sum_{x \in \text{misclassified}} \delta_x x$$

$p(t)$ - is decreasing ~~for~~ wrt (t) & $p(0) \leq 1$

Perceptron learning: An example:

Actually $x \in w_1$ But at t , $w(t)$ classify it as w_2

Let x be misclassified sample at time 't'. ~~ie w(t+1)~~

i.e $x \in w_1$ But classified by $w(t)$ as w_2 sample.

$\therefore \delta_x x = -x$ as $\delta_x = 1$ for w_2
 -1 for w_1

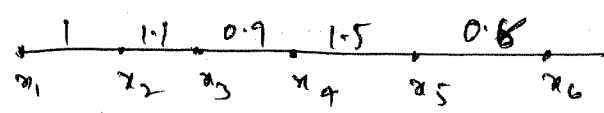
$\therefore w(t+1) = w(t) - P_t \sum_{x \in Y} \delta_x x$

Let $P_t = 1$ & x - is only misclassified sample.

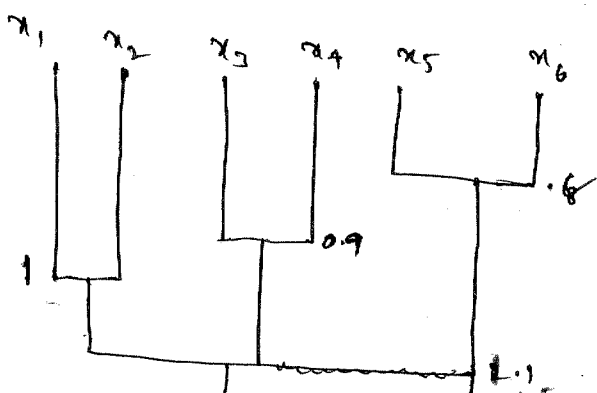
$w(t+1) = w(t) + x \Rightarrow$ By graphical addition direction of

~~ie~~ $w(t+1)$ & corresponding classification is shown.

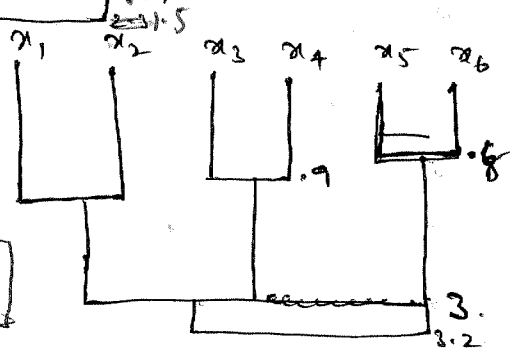
Example: Let $\{(0,0), (0,1), (-1,0)\} \in w_1$ and $\{(1,0), (0,1), (1,1)\} \in w_2$
 apply linear perceptron to evaluate w , with initial $w_0 = (-1,-1)^T$ & $b = 0$ & $P = 0.5$.



Single-link:



Complete link:



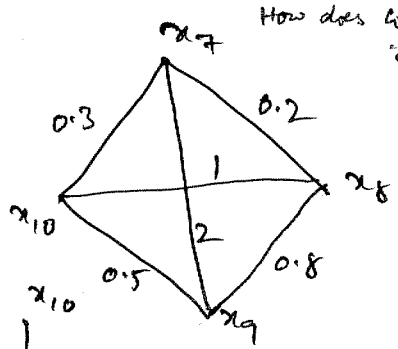
$D_{min}(D_i, D_j)$

$= \min_{x \in D_i, x' \in D_j} \|x - x'\|^2$

$D_{max}(D_i, D_j)$

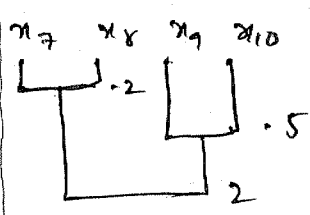
$= \max_{x \in D_i, x' \in D_j} \|x - x'\|^2$

$\{x_1, x_2, x_3, x_4\}, \{x_5, x_6\}$



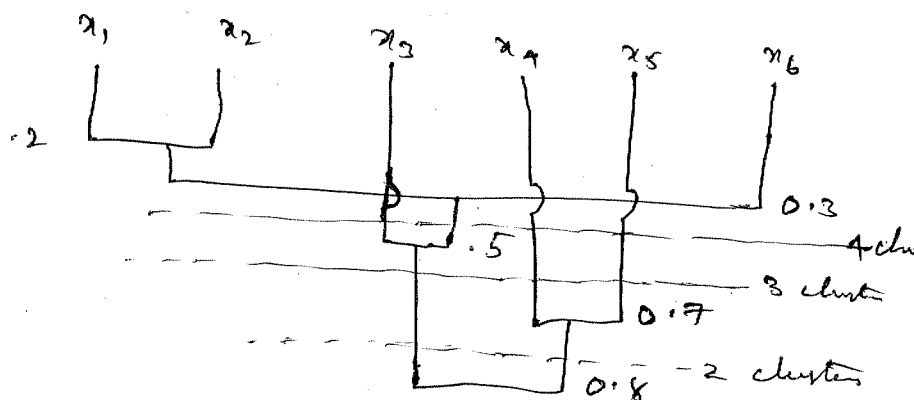
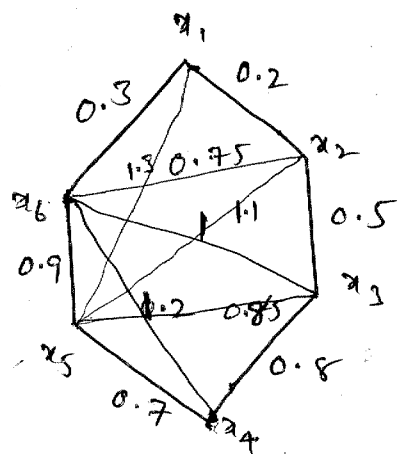
Breaking into two clusters:

$\{x_1, x_2, x_3, x_4\}, \{x_5, x_6\}$ or $\{x_7, x_8, x_{10}\}, x_9$



$\{x_7, x_8\}, \{x_9, x_{10}\}$

Single link:



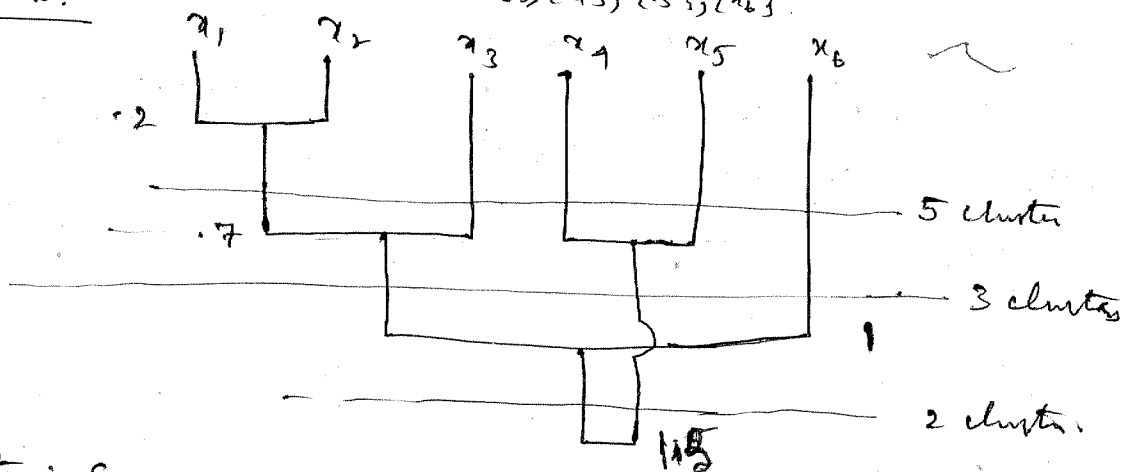
Two clusters: $\{x_1, x_2, x_3, x_6\}, \{x_4, x_5\}$

Three clusters: $\{x_1, x_2, x_3, x_6\}, \{x_4, x_5\}$

4-clusters: $\{x_1, x_2, x_6\}, \{x_3\}, \{x_4\}, \{x_5\}$

5-clusters: $\{x_1, x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_6\}$

Complete link:

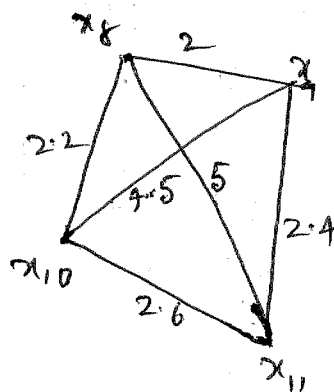
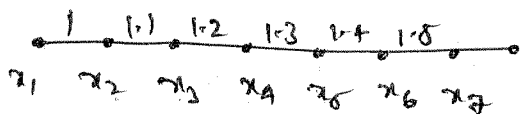


Two clusters: $\{x_1, x_2, x_3, x_6\}, \{x_4, x_5\}$

3 clusters: $\{x_1, x_2, x_3\}, \{x_4, x_5\}, \{x_6\}$

5 clusters: $\{x_1, x_2\}, \{x_3\}, \{x_4\}, \{x_5\}, \{x_6\}$

Problem:



→ Apply single link & complete link algorithms for clustering the points x_1 to x_{11} .

Introduction:

Neural networks:

The short coming of the linear 1 perceptron is that it can classify a data only if the samples of different classes are linearly separable.

* One of the approaches for classifying the nonseparable data is to apply kernel based methods i.e. non linear mappings of data in to a space where they are linearly

separable. ex:



if $x \in W_1$
 $0 \in W_2$

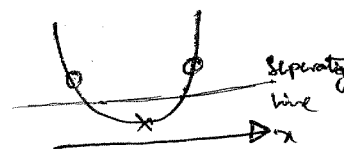
than data samples x_k are not linearly separable

we consider

But if $\phi(x_k) = e^{-\frac{x_k^2}{2\sigma^2}}$,

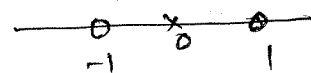
then

$\phi(x)$



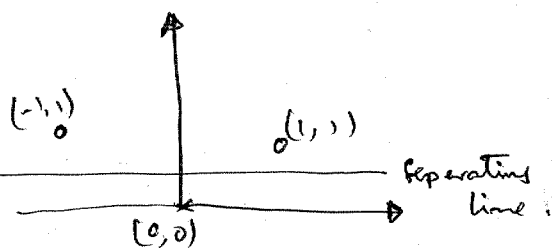
$\phi(x_k)$ are separable!

* Another approach is to increase the dimensionality so that the samples in the higher dimension are linearly separable. Considering same example x_k :



$x \in W_1$
 $0 \in W_2$

define $y = \begin{bmatrix} x \\ x^2 \end{bmatrix}$ then $\{y_k\} = \left\{ \begin{bmatrix} -1 \\ 1 \end{bmatrix}, \begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 \\ 1 \end{bmatrix} \right\}$



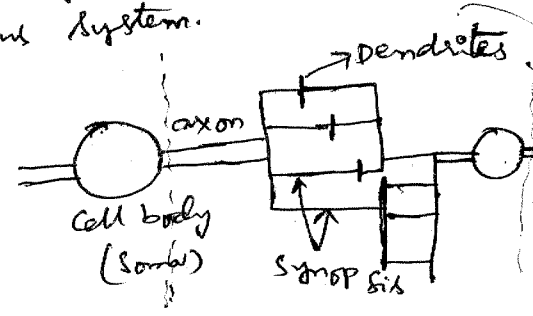
and $\{y_k\}$ are linearly separable.

However the choice of kernel function or arriving at higher dimensional representation ~~where~~ so that

the samples are linearly separable is not trivial to generalize and ~~depe~~ depends greatly on the data samples.

As against the techniques discussed so far, the neural networks based classification is motivated or inspired by ~~bio~~ Biological Capability of neurons (in the ~~human~~ animal brain). i.e. their nervous system.

Neurons have cell body, axon, Synapses & dendrites etc.



The spikes travelling along the axon triggers the release of neurotransmitter substance at the synapse of pre-synaptic neuron.

These neurotransmitters cause excitation or inhibition at the dendrites of post synaptic neuron.

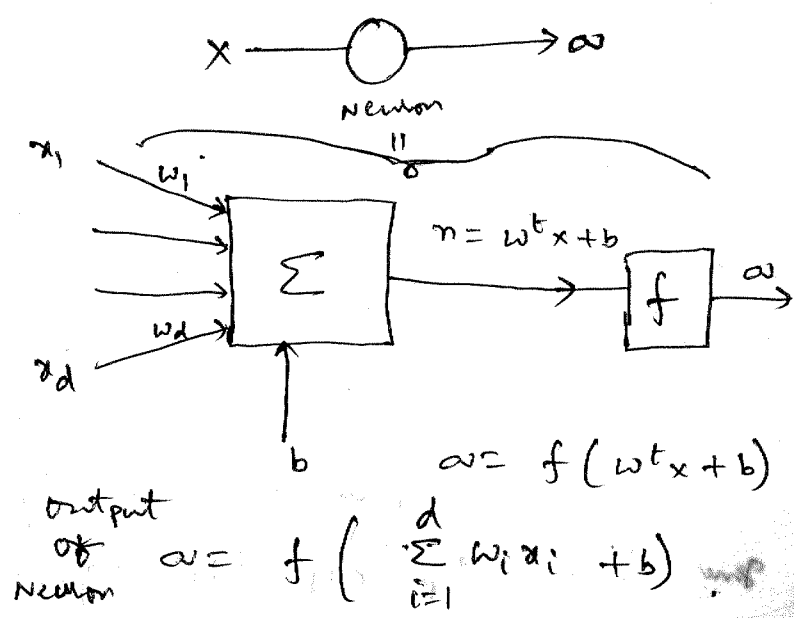
The integration of ^(weighted) excitations & inhibitions signals produce spikes in ^{axon of} post-synaptic neuron, which is modulated by ^{cell body by} transmitted to next neuron.

The contribution of signals depends on the strength of synaptic connections.

Now if we wish to develop a similar model

we note the following similarities:

neuron	its model
Strength of Synapse	weight
cell body	Transfer function and Summation
Signal in out put axon of neuron	output

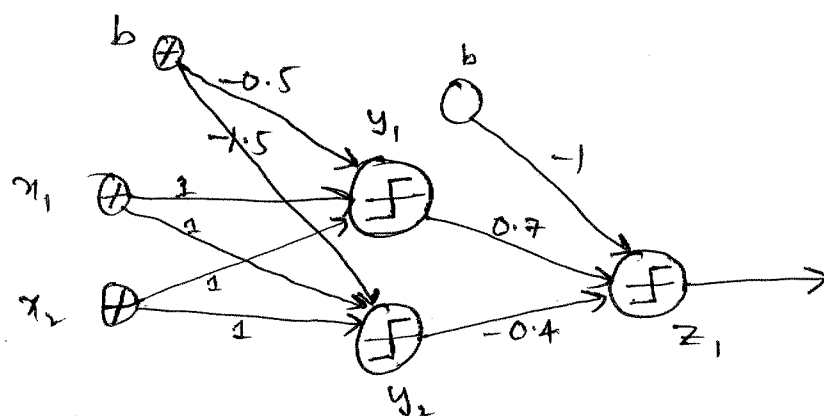


Feed Forward neural network:

(3)

* We consider following example to see the ability of neural network in classifying the non separable XOR data

~~In this~~

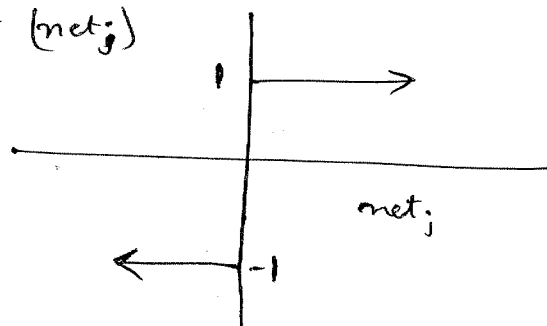


At hidden nodes

Let $net_j = \sum_{i=1}^d w_{ij} x_i \rightarrow$ weighted sum of inputs at j^{th} node/neuron

$y_j = f(net_j)$ Here f is given as threshold activation function.

$\Rightarrow f(net_j)$



i.e $f(net_j) = 1$ for $net_j \geq 0$
 -1 for $net_j < 0$

Similarly we defined at output node $net = \sum_{j=1}^2 y_j w_{jz}$
 $z = f(net)$; f again being threshold activation function. (here only one output & 2 hidden)

Note

$$net_1 = x_1 + x_2 - 0.5$$

$$net_2 = x_1 + x_2 - 1.5$$

$$\text{and } net = 0.7y_1 - 0.4y_2 - 1$$

Now XOR	x_1	x_2	net_1	y_1	net_2	y_2	net	z	NN. class. function
$w_2 \leftarrow$	0	0	-0.5	-1	-1.5	-1	-0.3	-1	w_2
w_1	0	1	.5	1	-0.5	-1	0.1	0	w_1
w_2	1	0	0.5	1	-0.5	-1	0.1	0	w_1
w_2	1	1	1.5	1	0.5	1	-0.7	-1	w_2

Feed forward neural network:

As from the above example

it is clear that as linear combination of perceptrons or an ~~set~~ organized set of neurons can separate the samples that are not linearly separable.

→ However, the question arise (i) how do we arrive at those weights that ^{optimally} separate the samples of two or more classes. And

→ (ii) what type of activation function to use at every neuron

(iii) How to choose the number of neurons their placement & number of layers etc.

Feed forward neural network has an input layer, one or more hidden layers & an output layer.

* no. of input layer neurons is proportional to no. of features

* no. of output layer neurons depends on no. of classes.

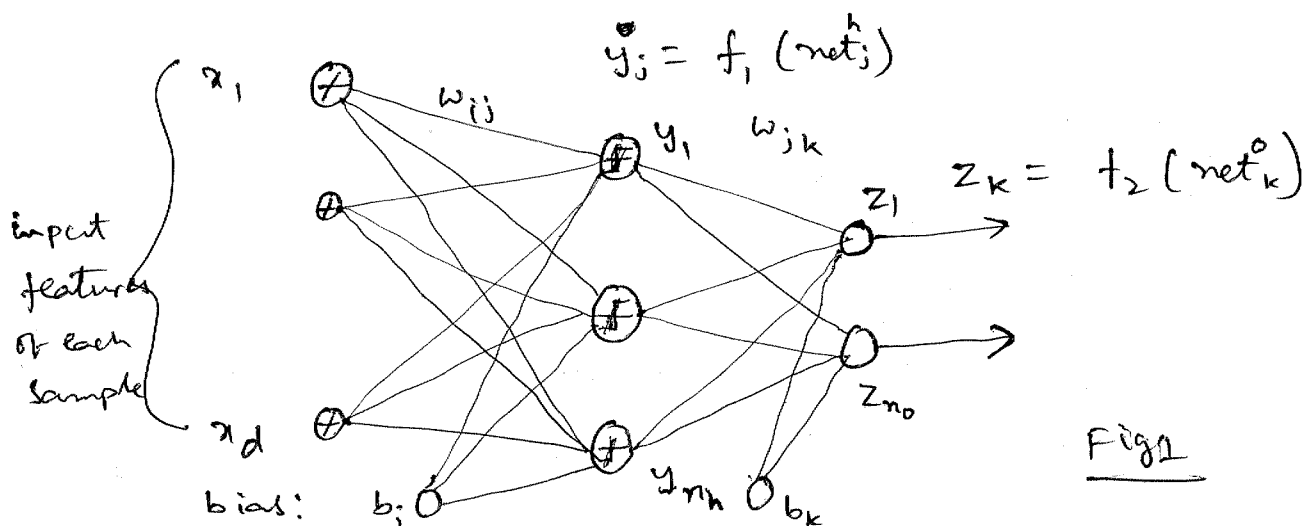


Fig 1

Single-layer hidden layer feed forward neural network.

looking at the neural network of Fig 1.
we can write the following:

$$\text{Eq 1} \left\{ \begin{array}{l} \text{net}_j^h = \sum_{i=1}^d x_i w_{ij} + b_j \\ y_j = f_1(\text{net}_j^h) \\ \text{net}_k^o = \sum_{j=1}^{n_h} y_j w_{jk} + b_k \quad ; \quad z_k = f_2(\text{net}_k^o). \end{array} \right.$$

where

- f_1 is the activation function at all neurons of ~~first~~ ^{hidden} layers.
- f_2 is the " " " " " of output layers.

eq 1 defines the mapping or the relation between input features & output value (or value), in terms of the synaptic weights & activation function (or biases).

i.e

$$z_k = f_2 \left[\sum_{j=1}^{n_h} w_{jk} f_1 \left(\sum_{i=1}^d w_{ij} x_i + b_j \right) + b_k \right]$$

→ now we address the problem of how to arrive at the weights that result in optimum classification.

→ There are several approaches for finding or learning the weights from training samples. The well known & widely used being the ^{error} Backpropagation algorithm.

→ As the name indicates we quantify the error at output neuron and back propagate it to update the weights by gradient descent.

(6)

Let d_k be the desired or target values at output node neurons for $(n)^{th}$ sample, then ~~total~~ mean square error at output layer of NN is (for given set of weights)

$$J(w) = \frac{1}{2} \sum_{k=1}^{n_0} (d_k - z_k)^2 \quad \text{or} \quad \frac{1}{2} \|d - z\|^2$$

factor $\frac{1}{2}$ is introduced for later simplifications (for convenience).

→ The backpropagation learning rule is based on gradient descent. The weights are ~~first~~ initialized with random values, and then changed in the direction that reduce the error (i.e. opposite to direction of error gradient).

$$\Delta w = -\eta \cdot \frac{\partial J}{\partial w} \quad \text{or} \quad \Delta w_{pq} = -\eta \frac{\partial J}{\partial w_{pq}}$$

η - is learning rate ($= \rho$). (in component form)

Then the weights are iteratively updated as

~~$$w_{k+1}$$~~
$$w(t+1) = w(t) + \Delta w(t).$$

→ To obtain the Δw we need to find $\frac{\partial J}{\partial w}$ w.r.t w_{ij} & w_{jk} .

→ Starting from output layer,

$$J = \frac{1}{2} \sum_{k=1}^{n_0} [d_k - f_2(\text{net}_k)]^2$$

we need $\frac{\partial J}{\partial w_{kj}} = - \frac{\partial J}{\partial \text{net}_k} \frac{\partial \text{net}_k}{\partial w_{kj}} = [d_k - f_2(\text{net}_k)] f_2'(\text{net}_k).$

$$\text{net}_k = \sum_{j=1}^{n_h} w_{kj} y_j + b_j, \quad \frac{\partial J}{\partial \text{net}_k} =$$

→ Starting from output layer: at n^{th} or t^{th} iteration: (7)

we have $J(A) = \frac{1}{2} \sum_{k=1}^{n_0} \left[d_k - f_2(\text{net}_k) \right]^2$

we can write

$e(n) = d_k - f_2(\text{net}_k)$ for finding Δw_{jk}

we find $\frac{\partial J}{\partial w_{jk}} = \frac{\partial J}{\partial \text{net}_k^o} \cdot \frac{\partial \text{net}_k^o}{\partial w_{jk}} = -\delta_k \frac{\partial \text{net}_k^o}{\partial w_{jk}}$

$\text{net}_k^o = \sum_{j=1}^{n_h} w_{jk} y_j \Rightarrow \frac{\partial \text{net}_k^o}{\partial w_{jk}} = y_j$

$$-\delta_k^o = \frac{\partial J}{\partial \text{net}_k^o} = \frac{\partial}{\partial \text{net}_k^o} \left[\frac{1}{2} \sum_{k=1}^{n_0} \left[d_k - f_2(\text{net}_k) \right]^2 \right]$$

$$= - \left[d_k - f_2(\text{net}_k) \right] f_2'(\text{net}_k) = -(d_k - z_k) f'(\text{net}_k)$$

 - (1)

$\Rightarrow \frac{\partial J}{\partial w_{jk}} = - (d_k - z_k) \cdot f'(\text{net}_k) \cdot y_j$

$\Delta w_{jk} = \cancel{\eta e_k} \rightarrow \eta e_k \cdot f'(\text{net}_k) \cdot y_j$

→ Now for updating the weights of inputs to hidden layer neurons, we need $\Delta w_{ij} = -\eta \frac{\partial J}{\partial w_{ij}}$

$$\frac{\partial J}{\partial w_{ij}} = \frac{\partial}{\partial w_{ij}} \cdot \frac{1}{2} \sum_{k=1}^{n_0} \left[d_k - f_2 \left(\sum_{j=1}^n w_{jk} \frac{f_1(\text{net}_j^h)}{y_j + b_j} \right) \right]^2$$

$$= \frac{\partial J}{\partial \text{net}_j^h} \cdot \frac{\partial \text{net}_j^h}{\partial w_{ij}} = -\delta_j^h \frac{\partial \text{net}_j^h}{\partial w_{ij}} \quad \text{--- (2)}$$

$$\frac{\partial J}{\partial \text{net}_j^h} = \frac{\partial J}{\partial y_j} \cdot \frac{\partial y_j}{\partial \text{net}_j^h}; \quad \frac{\partial J}{\partial y_j} = \sum_{k=1}^{n_0} (d_k - z_k) \cdot f_2'(\text{net}_k) \cdot w_{jk}$$

$$\frac{\partial y_j}{\partial \text{net}_j^h} = f_1'(\text{net}_j^h) \quad \text{--- (4)} \quad = - \sum_{k=1}^{n_0} \delta_k w_{jk} \frac{\partial J}{\partial y_j} \quad \text{from (1)} \rightarrow \text{--- (3)}$$

$$\frac{\partial \text{net}_j}{\partial w_{ij}} = x_i \quad - (5)$$

(8)

from (1), (3), (4) & (5), we have

$$\frac{\partial J}{\partial w_{ij}} = - \underbrace{\sum_{k=1}^{n_0} \delta_k^o w_{ik}}_{\delta_j^h} f'_j(\text{net}_j^h) \cdot x_i$$

$$\begin{aligned} \therefore \Delta w_{ij} &= -\eta \frac{\partial J}{\partial w_{ij}} = \eta \delta_j^h x_i \\ &= \eta \left[\sum_{k=1}^{n_0} w_{ik} \delta_k^o \right] \cdot f'_j(\text{net}_j^h) \cdot x_i \end{aligned}$$

Thus based on these rules we can update the weights at different layers.

→ Now comes to the selection of activation functions

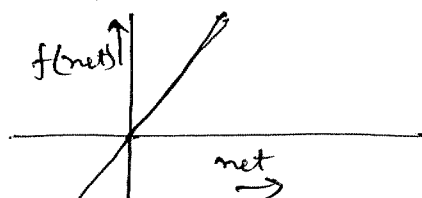
'f' (or f_1 & f_2) as required in weight update, it must be differentiable.

(ii) to bring proportional change in the output, from its inputs, it must be monotonically increasing.

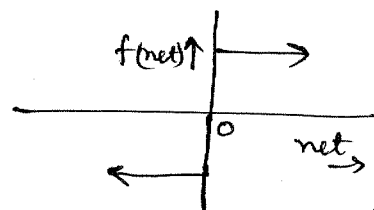
(iii) it must not be too high or too low → well within range to take a decision

A few common and widely used activation functions:

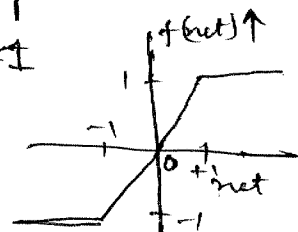
(i) linear: $f(\text{net}) = \text{net}$



(ii) Threshold $f(\text{net}) = \begin{cases} 1 & \text{for } \text{net} \geq 0 \\ 0 & \text{for } \text{net} < 0 \end{cases}$



(iii) Piecewise linear: $f(\text{net}) = \begin{cases} \text{net} & \text{for } -1 < \text{net} < 1 \\ -1 & \text{for } \text{net} < -1 \\ 1 & \text{for } \text{net} > 1 \end{cases}$



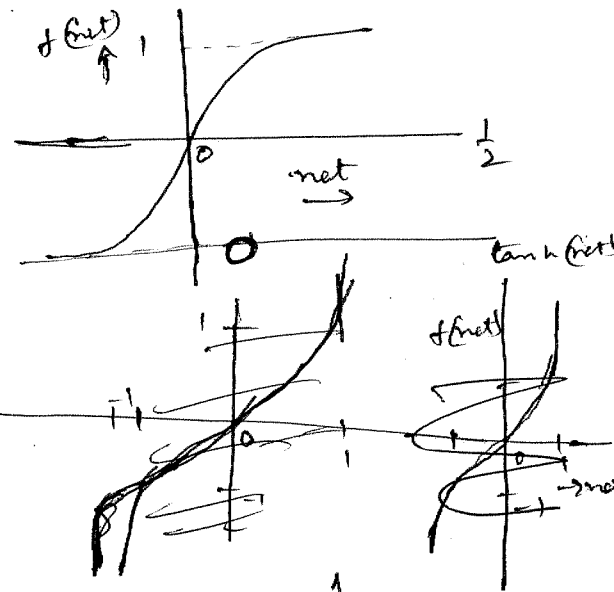
Activation function:

(iv) Sigmoid function: $f(\text{net}) = \frac{1}{1 + e^{-\text{net}}}$ as $\text{net} \rightarrow \infty$
or also called logistic function.

(v) tan hyperbolic function:

$$f(\text{net}) = \frac{e^{\text{net}} - e^{-\text{net}}}{e^{\text{net}} + e^{-\text{net}}}$$

$$\tanh(\text{net}) = \frac{e^{2\text{net}} - 1}{e^{2\text{net}} + 1}$$



Choice of no. of hidden neurons:

- Higher the ~~more~~ complexity of decision boundaries, larger the no. of hidden neurons required
- As we have discussed & if the samples (of two classes) are linearly separable single perceptron or (neuron with threshold activation function) is sufficient to classify correctly.
- However as the samples become more and more intermingled or not separable linearly, larger no. ~~new~~ hidden neurons are needed.

Choice of input & output neurons: Choice of input neurons

usually depends on the dimensionality of features of each sample. if d -features are attributed with each sample d -input neurons are used. ~~the~~ while the no. of output neurons depends on the no. of classes / clusters in the data. If we

we sigmoid output neuron more ~~no. of~~ ~~input~~ even with one or few output neuron more classes can be mapped.

However if we use a threshold activation function with output neurons, then each output neuron can classify only two classes.

no. of hidden layers: depends on the ^{problem} domain, if we want to

~~Settling in local minima~~ incorporate model operations like Translation

Heuristic Rotation & Scaling each can be modeled by one-layer.

Weight update: Even if we use large network with large no. of neurons, the weight update

$$w(m) = (1 - \epsilon) w(m), \text{ where } \epsilon \text{ is } 0 < \epsilon < 1$$

helps to nullify the smaller, possibly unimportant weights

Choice of no. of hidden layers: usually single

hidden layer NN with arbitrarily large no. of neurons must be sufficient to model any nonlinearity in mapping or classification. However if the problem/task at hand demands few independent or various operation like translation, rotation & scaling to be incorporated through training processes, then the multiple hidden layers NN is preferred. Each layer can be trained to effectively model each operation.

overfitting: If we increase the no. of hidden neurons or

layers to a very large extent then there is a possibility of getting the decision boundaries/mappings more tuned

only to the training samples. In that case we may have good ~~best~~ performance with training set, but fails and does not

provide ^{even} minimum performance on the test set. The decision boundaries will be more local than needed. This to some extent can be overcome by ~~using~~ ~~the~~ ~~same~~ ~~data~~ ~~sets~~ ~~for~~ ~~training~~ ~~and~~ ~~testing~~ with cross validation ~~for~~ ~~too~~ and continuing to minimize errors that give fairly good performance.

→ Heuristic weight update: will also help in avoiding over fitting. This eliminates the role of ^{ineffective} neurons (or makes them dummy neurons) by decreasing their weight (initially) to almost zero values. This is

done by $w(m) = \epsilon w(m) \quad 0 < \epsilon \leq 1$
 \uparrow near to 1

By this smaller weights gets smaller & smaller & subsequently eliminated. while as effective weights pick up well in gradient descent update process.

Local minima vs global minima: If the cost function has

several local minima then the net may end up ~~with~~ by reaching local minima, based on how the

weights are initialized (and how they are updated).

Momentum term: In error surface there can be series where $\frac{dJ}{dw} \approx 0$ to overcome ~~local minima~~ ^{these} a momentum term may be

(Plateaus)

added as ~~$w(m+1) = w(m) + \Delta w(m)$~~

$$w(m+1) = w(m) + \Delta w(m) + \alpha \Delta w(m-1)$$

$\Delta w(m)$ - is gradient descent update

while $\Delta w(m-1) = w(m) - w(m-1)$.

Then the movement of weights will be more smoother.

$$0 < \alpha < 1.$$

B. Standardizes input features: Before feeding into the neural network, the input features are to be ~~the~~ scaled & demeaned so as to ensure that all features are in same range. If they are not the features at higher dynamic range / scale will be given more contribution in update as
$$\left(\text{net} = \sum_{i=1}^d w_{ij} x_i \right)$$
, ~~and~~ weights are uniformly initialized. So effectively we may need to apply whitening transform i.e. make ^{all} features zero mean & unit variance, before feeding to NN.

Initializing weights: ~~but~~ weights are usually randomly initialized. (from backpropagation update eq. it is clear that if $w_{ij} = 0$ then $\Delta w_{ij} = 0$ and update does not take place. So they can't be zero). usually the weights are uniformly distributed in such a way that $-1 < \underbrace{\sum_{i=1}^d w_{ij} x_i}_{\text{net}} < 1$ as activation function become linear ⁱⁿ ~~at~~ this range. if weights are high than $|\text{net}| > 1$ & saturation occurs. Thus we prefer to ~~be~~ initialize the weights in such a way that $\left| \sum_{i=1}^d w_{ij} x_i \right| < 1$. if we have d -dimensions it can be shown that

~~Set~~ If $-\frac{1}{\sqrt{d}} < w_{ij} < \frac{1}{\sqrt{d}}$ $\forall i, j$ then

$$-1 \leq \sum_{i=1}^d w_{ij} x_j \leq 1$$

Thus we initialize all weights, by uniform distribution basis in range $\left[-\frac{1}{\sqrt{d}}, \frac{1}{\sqrt{d}}\right]$.

for hidden to output neurons we initialize in same $\left[-\frac{1}{\sqrt{n_h}}, \frac{1}{\sqrt{n_h}}\right]$.

Non uniform convergence & Bayes Rate:

If the weights are nonuniformly initialized there is a possibility that few ~~fewer~~ ^{fewer} reach their optimal value faster than others & non uniform convergence of NN occurs, this makes the output error to differ. marked by from Bayes error rate. we use the above discussed ~~into~~ ^{uniform} distribution based initialization to overcome this.

Training protocols: There are three broad categories of

Training NN. They are stochastic, batch and on line learning

In stochastic training we present ~~in~~ training patterns randomly from training set and the network weights are updated for each

pattern. In batch ~~train~~ training we present all the patterns ^{corresponding weight ~~are~~ ^{are} summed up then only} are presented to network once, ~~each~~ each weight update takes place.

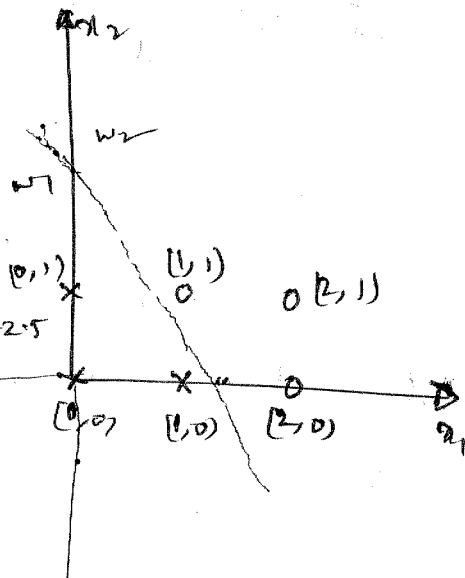
In both of these trainings we must virtually make several passes through the training set each pass is called an epoch. In online training each pattern is presented once & only once.

4.

$$W = \begin{bmatrix} -1 \\ -1 \end{bmatrix} \quad W_0 = 0.5, \quad P=1$$

$$W^t x = \begin{matrix} .5 & -.5 & -.5 & -2.5 & -3.5 & -2.5 \end{matrix}$$

misclassified



$$W = \begin{bmatrix} -1 \\ -1 \\ 0.5 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 2 \end{bmatrix} = \begin{bmatrix} 0 \\ 0 \\ 2.5 \end{bmatrix}$$

$$W = \begin{bmatrix} 0 \\ 0 \\ 2.5 \end{bmatrix} - \begin{bmatrix} 5 \\ 2 \\ 3 \end{bmatrix} = \begin{bmatrix} -5 \\ -2 \\ -5 \end{bmatrix}$$

-5 -2.5 -5.5

$$W = \begin{bmatrix} -5 \\ -2 \\ -5 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \\ 3 \end{bmatrix} = \begin{bmatrix} -4 \\ -1 \\ 2.5 \end{bmatrix}$$

2.5 -1.5 1.5 -6.5 -ve -ve

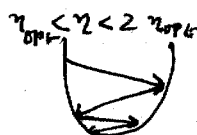
$$W = \begin{bmatrix} -4 \\ -1 \\ 2.5 \end{bmatrix} + \begin{bmatrix} 1 \\ 0 \\ 1 \end{bmatrix} = \begin{bmatrix} -3 \\ -1 \\ 3.5 \end{bmatrix}$$

$$W^t x + W_0 = 0 \Rightarrow -3x_1 - x_2 + 3.5 = 0$$

$$\Rightarrow 3x_1 + x_2 - 3.5 = 0 \quad x_2 = 0 \Rightarrow x_1 = \frac{3.5}{3}, \quad x_1 = 0 \Rightarrow x_2 = 3.5$$

Stopping criterion; $\|\nabla J(W)\| < \theta$, where $J = \sum_{p=1}^n J_p$

optimal learning rate: $\eta_{opt} = \left(\frac{\partial J}{\partial \eta} \right)^{-1}$



Converge for $\eta < 2\eta_{opt}$

Slow for $\eta < \eta_{opt}$ or $\eta_{opt} < \eta < 2\eta_{opt}$

Diverge for $\eta > 2\eta_{opt}$

activation function's prop: Saturation \rightarrow outputs mapped to posteriors

Polynomial activation, radial basis/gaussian, monotonicity \rightarrow
sigmoid activation, $\frac{e^{b_{net}} - e^{-b_{net}}}{e^{b_{net}} + e^{-b_{net}}}$ nonlinear, differentiable \rightarrow imp.
parameter selection \rightarrow generating/creating manufacturing training data: Adds d-dimensional Gaussian noise

ex by rotation, scale, etc in character recognition.

expressive power: no. of weights

Pattern Recognition

VI – Physical science

Assignment – 5

18-4-2012

1. Let $p(\mathbf{x}|\omega_i)$ be arbitrary densities with means μ_i and covariance matrices Σ_i — not necessarily normal — for $i = 1, 2$. Let $y = \mathbf{w}'\mathbf{x}$ be a projection, and let the induced one-dimensional densities $p(y|\omega_i)$ have means μ_i and variances σ_i^2 .
 - (a) Show that criterion function $J_1(\mathbf{w}) = (\mu_1 - \mu_2)^2 / (\sigma_1^2 + \sigma_2^2)$ is maximized by $\mathbf{w} = (\Sigma_1 + \Sigma_2)^{-1}(\mu_1 - \mu_2)$.
 - (b) If $P(\omega_i)$ is the prior probability for ω_i , show that $J_2(\mathbf{w}) = (\mu_1 - \mu_2)^2 / (P(\omega_1)\sigma_1^2 + P(\omega_2)\sigma_2^2)$, is maximized with $\mathbf{w} = [P(\omega_1)\Sigma_1 + P(\omega_2)\Sigma_2]^{-1}(\mu_1 - \mu_2)$.
2. Given samples $x_1 = [0,0]$, $x_2 = [1,0]$, $x_3 = [0,3]$, $x_4 = [-2,0]$, $x_5 = [-3,0]$, $x_6 = [-3,3]$. Find the dissimilarity matrix, (i) using Euclidian distance (ii) City block distance (defined as $d(x_1, x_2) = |dx| + |dy|$, where dx, dy are distances in x and y , directions respectively), and find the similarity matrix using (iii) cosine angle or correlation measure and (iv) tonimato measure.
3. Apply single link and complete link clustering algorithms to samples in above problem based on the (i) Euclidian distance based dissimilarity matrix and (ii) Tanimoto measure based similarity matrix. Break in to clusters that contain less than or equal to 3 elements in each cluster. Comment on the kind of clusters resulted by these methods.
4. Given samples $x_1 = [0,0]$, $x_2 = [1,0]$, $x_3 = [0,1]$, $x_4 = [2,0]$, $x_5 = [2,1]$, $x_6 = [1,1]$ where first 3 samples belongs to class1 and next 3 samples to class2, and initial direction of the linear discriminant function $w_0 = [-1, -1]$. Find the direction of linear discriminant function that separates the class 1 and 2 samples, by perception learning, (i) with learning rate of 1 and 0.5 (assume the samples on the separating hyper-plane as correctly classified). (ii) Repeat learning or find decision hyper-plane, with a learning rate of 0.5, if the samples on the separating hyper-plane are assumed to be mis-classified.
5. Represent a single neuron as perception and discuss its limitations. Draw a neural network that can be used for XoR classification task.
6. Take a 2-2-1 neural network, take same weights and bias at all neurons as used in class, but use sigmoid activation function at all hidden neurons and threshold activation function at output neuron, and verify whether it solves XoR problem or not?
7. Represent a feed forward neural network with d -input, n -hidden and c -output neurons, designate connection weights, activation functions and outputs. Assume you have a training set with known desired outputs d_k , $k = 1:c$. Derive the weight update equations by back propagation algorithm.
8. Comment on the choice of number of input neurons, output neurons and hidden neurons. How do you initialize weights and if the input features ^{are} at different scales/have different dynamic range how will you present them to NN? What is the problem of over fitting in NN and how will you minimize it? How do you asses performance of your NN?
9. What is the importance of momentum term and how do we chose activation function. Give few typical choices. What are different types of training NNs.

take
bias
 $w_0 = 0.5$

Support Vector Machine / Maximum margin classifier:

①

linear discriminant function: $y(x) = w^T x + w_0$

x assigned to class 1 if $y(x) \geq 0$
else assigned to class 2.

\therefore decision boundary is $y(x) = 0$ or $w^T x + w_0 = 0$. It can be line, plane, or hyper plane.

It can be seen that for two points on hyper plane, $y(x_A) = 0$,

$y(x_B) = 0 \Rightarrow w^T (x_A - x_B) = 0$ & w is \perp to every line on decision plane.

Similarly if ' x_0 ' is a point on decision surface, then the normal distance of decision surface from origin is

given by

$$\frac{w^T x_0}{\|w\|} = -\frac{w_0}{\|w\|}$$

Since $y(x_0) = 0$. Thus the bias parameter determines the location of decision plane.

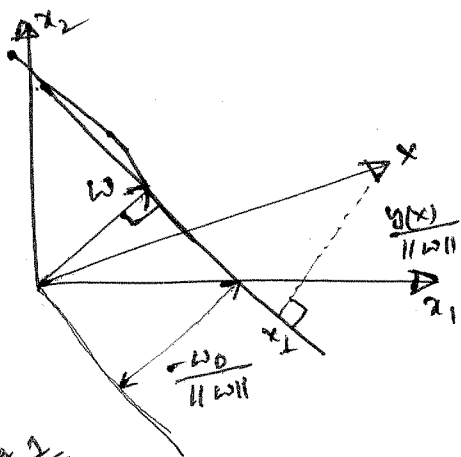


Fig 2

\rightarrow if ' x ' is not on decision surface, then $\frac{y(x)}{\|w\|}$ gives

the signed \perp distance of point from decision plane.

This can be proved as follows.

if $x = x_1 + \gamma \frac{w}{\|w\|}$, where x_1 is on decision plane.

$y(x) = y(x_1) + \gamma$ multiplying both sides by w & ~~and~~

adding w_0 we get $w^T x + w_0 = \underbrace{w^T x_1 + w_0}_{y(x_1)} + \gamma \frac{w^T w}{\|w\|}$

$$\therefore \gamma = \frac{y(x)}{\|w\|}$$

Maximum margin classifier:

(2)

$$y(x) = w^T \phi(x) + b$$

We decide w_1 if $y(x) > 0$, w_2 otherwise.

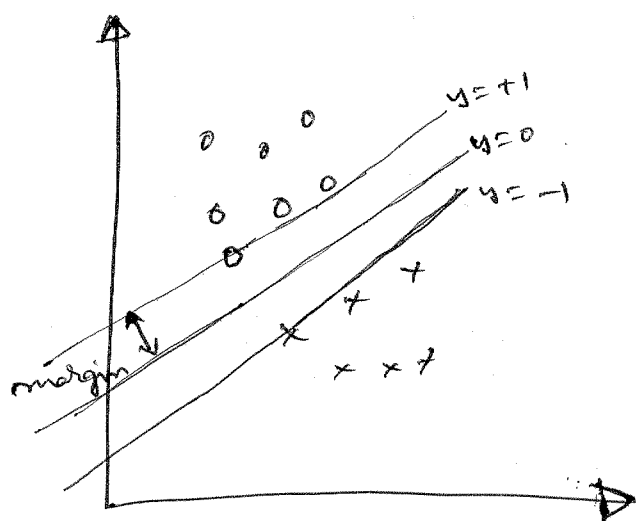
Let the training data set comprises of x_1, \dots, x_n , with labels or corresponding target values t_1, \dots, t_n , where $t_n \in \{-1, 1\}$

such that $t_n = 1$ for $y(x_n) > 0$

$t_n = -1$ for $y(x_n) < 0$

assuming all the points are linearly separable, i.e. there exist at least one w & b to separate class 1 & 2 data samples.

observe that $t_n y(x_n) \geq 0$ ~~always~~ for all data points.



In general several solutions can exist for w , that separates the two class data. Here we are interested in such a solution that maximize the 'margin' w.r. to two class data points. Here margin is defined as the \perp er

distance from closest point (of both classes) to the plane separating the two class data. In Fig 1

We quantified the ^{signed} \perp er distance of any point to hyper plane

$y(x)=0$ as $\frac{y(x)}{\|w\|}$. Thus the (absolute) distance

of point x_n to the decision surface is given by

$$\frac{t_n y(x_n)}{\|w\|} = \frac{t_n (w^T \phi(x_n) + w_0)}{\|w\|}$$

This margin ~~$t_n (w^t \phi(x_n) + b)$~~

$$\frac{t_n (w^t \phi(x_n) + b)}{\|w\|}$$

gives the \pm distance from plane to closest point x_n in the data set, we wish to optimize ~~the~~ ^{wrt} parameters 'w' & 'b' in order to maximize the margin. The maximum margin solution is found by solving

$$\arg \max_{w, b} \left\{ \frac{1}{\|w\|} \cdot \min_n (t_n w^t \phi(x_n) + b) \right\}$$

maximizing margin. \uparrow \uparrow wrt closest point x_n

if we set $w \rightarrow kw$ & $b \rightarrow kb$ then the distance with x_n to decision surface $\frac{y(x_n)}{\|w\|}$ will not change thus we utilize this freedom to set

$$t_n (w^t \phi(x_n) + b) = 1 \quad - (1)$$

for the closest point. For all other points

$$t_n (w^t \phi(x_n) + b) \geq 1 \quad - (2)$$

The constraint is had to be active, if the equality holds. There will always be one such closest point. And when we maximize the margin (wrt to two class data points) there will at least be two data points ~~for~~ (one from each class) for which the constraint

(2) is active. on the other hand we need to ~~minimize~~ ^{maximize} $\|w\|^2$. so we have to solve the optimization problem

$$\arg \min_{w, b} \frac{1}{2} \|w\|^2$$

Constraint (2) ^{factor} $\frac{1}{2}$ is for mathematical convenience.

To do this we construct a Lagrange multiplier function (4)

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{n=1}^N \alpha_n \left\{ t_n (w^T \phi(x_n) + b) - 1 \right\}$$

Where $\alpha = (\alpha_1, \dots, \alpha_N)^T$. (4)

-ve sign in front of second term is because, though we need to minimize wrt w & b , we are maximizing wrt α_n . Setting derivation of $L(w, b, \alpha)$ wrt w & b to

zero we get

$$w = \sum_{n=1}^N \alpha_n t_n \phi(x_n) \quad -5$$

$$\text{and } \sum_{n=1}^N \alpha_n t_n = 0 \quad -6$$

Using (5) & (6) & eliminating w & b from $L(w, b, \alpha)$ gives dual representation of the maximum margin problem. ~~In which~~

$$\tilde{L}(w, b, \alpha) = \frac{1}{2} \langle w, w \rangle - \sum_{n=1}^N \alpha_n t_n (w^T \phi(x_n) + b) + \sum_{n=1}^N \alpha_n$$

$$\begin{aligned} \tilde{L}(\alpha) &= \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \langle \phi(x_n) \cdot \phi(x_m) \rangle \\ &\quad - \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \langle \phi(x_n) \cdot \phi(x_m) \rangle + 0 + \sum_{n=1}^N \alpha_n \end{aligned}$$

$$\tilde{L}(\alpha) = \sum_{n=1}^N \alpha_n - \frac{1}{2} \sum_{n=1}^N \sum_{m=1}^N \alpha_n \alpha_m t_n t_m \overset{k(x_n, x_m)}{\langle \phi(x_n) \cdot \phi(x_m) \rangle}$$

$$\text{where } k(x_n, x_m) = \langle \phi(x_n) \cdot \phi(x_m) \rangle = \phi(x_n)^T \phi(x_m)$$

This dual we need to maximize with the constraints

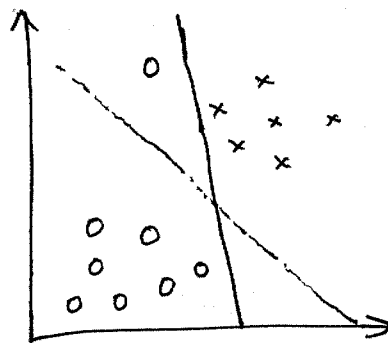
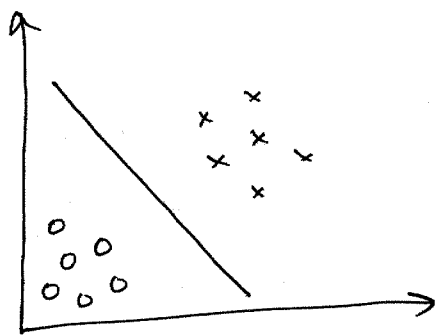
$$\alpha_n \geq 0 \quad \forall n=1, \dots, N \quad \text{and} \quad \sum_{n=1}^N \alpha_n t_n = 0.$$

In order to classify, new data point, with trained model,

$$\text{we have } y(x) = w^T x + b = \sum_{n=1}^N \alpha_n t_n k(x, x_n) + b.$$

Support vector Machine: Regularization in non Separable case! ①

While mapping data to high dimensional feature space via ϕ does generally increase the likelihood that the data is separable, but does not guarantee that it always will be so. Also, in some cases it is not clear that finding a separating hyperplane is exactly what we would want to do, since that might be susceptible to outliers. For instance, the left figure below shows an optimal margin classifier, and when a single outlier is added in the upper left region (right figure), it causes the decision boundary to make a dramatic swing, and the resulting classifier has a much lesser margin.



To make the algorithm work for non-linearly separable data sets as well as be less sensitive to outliers, we reformulate our optimization (using l_1 regularization) as follows:

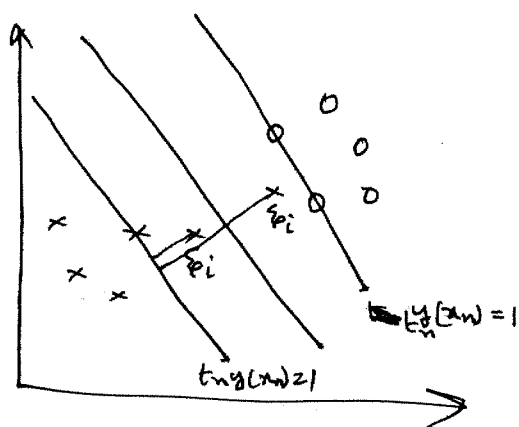
$$\min_{w, b} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^m \xi_i$$

$$\text{such that } w^t \phi(x_n) + b \geq 1 - \xi_n, \quad n=1, \dots, m.$$
$$\xi_n \geq 0, \quad n=1, \dots, m.$$

Thus examples are now permitted to have functional margin less than 1. And if an example has a functional margin $1 - \xi_n$

(2)

(with $\xi > 0$), we would pay a cost of the objective function being increased by $C \xi_i$.



if for a sample $t_n y(x_n) = 1 - \xi_i$, $\xi_i > 0$

$C \xi_i \rightarrow$ is its associated cost.

if $0 < \xi_i < 1$ it is correctly classified but in between the ^{separating} hyperplane & margin. If $\xi_i > 1$ then it is an outlier, wrongly classified with associated cost $C \xi_i$.

The Parameter C controls the relative weighting between the twin goals of making the $\|w\|^2$ small (which makes the margin large) and of ensuring that most examples have functional margin at least 1.

For this non-separable case, we can form the Lagrangian:

$$L(w, b, \xi, \alpha, \gamma) = \frac{1}{2} w^T w + C \sum_{i=1}^m \xi_i - \sum_{i=1}^m \alpha_i [t_i (\phi(x)^T w + b) - 1 + \xi_i] - \sum_{i=1}^m \gamma_i \xi_i \quad \text{--- (D)}$$

Here, the α_i 's & γ_i 's are Lagrange multipliers (constrained to be ≥ 0). After setting the derivatives wrt ξ_i, w & b to zero as before, substituting them back in, & simplifying, we get following dual form of the problem:

$$\max_w L(w) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m t_i t_j \alpha_i \alpha_j \langle \phi(x_i), \phi(x_j) \rangle$$

such that $0 \leq \alpha_i \leq C, \quad i=1, \dots, m$

$$\sum_{i=1}^m \alpha_i t_i = 0$$

as before $w = \sum_{i=1}^m \alpha_i t_i \phi(x_i)$ and $y(x) = w^T x + b$

$$= \sum_{i=1}^m \alpha_i t_i \langle x_i, x \rangle + b$$

$$y(x) = w^T x + b$$

$$= \sum_{i=1}^m w_i t_i \langle x_i, x \rangle + b.$$

(3)

We note that, somewhat surprisingly, in adding L₂ regularization the only change to the dual problem is that what was originally constrained that $0 \leq w_i$ has now become $0 \leq w_i \leq C$.

This dual problem can be solved by sequential minimal optimization (SMO).

Sequential ~~optimal~~ minimal optimization (SMO): The SMO algorithm gives an efficient way of solving dual problems arising from the derivation of SVM. It is based on coordinate ascent

optimization algorithm.

Coordinate ascent: Consider trying to solve the unconstrained

optimization problem $\max_{\alpha} W(\alpha_1, \alpha_2, \dots, \alpha_m)$

Here, we think of W as just some function of the parameters α_i 's. and for now we ignore any relation ship between this problem & SVMs. The coordinate ascent can be summarized as

```

loop until convergence {
    for  $i = 1, \dots, m$  {
         $\alpha_i : \arg \max_{\hat{\alpha}_i} W(\alpha_1, \dots, \alpha_{i-1}, \hat{\alpha}_i, \alpha_{i+1}, \dots, \alpha_m).$ 
    }
}

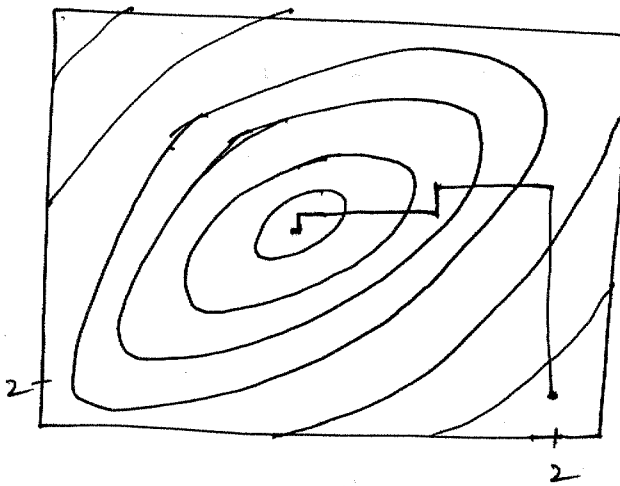
```

Thus, in the innermost loop of this algorithm, we will hold all the variables except for some α_i fixed, and re-optimize W with respect to just the parameter α_i . Similarly each at a time algorithm re-optimizes w.r.t $\alpha_1, \dots, \alpha_m$, holding all others fixed.

(4)

A more sophisticated version might choose the ordering. For instance, we may choose the next variable to update according to which one we expect to allow us to make the largest increase in $w(\alpha)$.

When the function w happens to be of such a form that the "arg max" in the inner loop can be performed efficiently then coordinate ascent can be a fairly efficient algorithm. Here is a picture of coordinate ascent in action:



The ellipses in the figure are the contours of a quadratic function that we want to optimize. Coordinate ascent is initialized at $(2, -2)$ and also plotted in

figure is the path it took to on its way to the global maximum. We note that coordinate ascent takes a step that's parallel to one of the axes, since only one variable is optimized at a time.

SMD: Now we give details of SMD algorithm to solve dual problem in SVMs. Here is the dual optimization problem that we want to solve. (assume $N=m$)

$$\max_{\alpha} w(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j t_i t_j \langle \underline{x_i}, \underline{x_j} \rangle \quad -①$$

$$\text{s.t.} \quad 0 \leq \alpha_i \leq C, \quad i=1, \dots, N \quad -②$$

$$\text{and} \quad \sum_{i=1}^N \alpha_i t_i = 0 \quad -③$$

Let's say that we have ω_i 's satisfying ② & ③. Now suppose we want to hold ~~ω_1~~ $\omega_2, \dots, \omega_N$ fixed, and take a coordinate ascent step and re-optimize the objective with respect to ω_1 . Can we make any progress? The answer is no, because the constraint ③ ensures

$$\sum_{i=1}^N \omega_i t_i = 0 \Rightarrow \omega_1 t_1 = - \sum_{i=2}^N \omega_i t_i$$

$$\Rightarrow t_1 = - \sum_{i=2}^N \omega_i t_i$$

(since $\omega_i^2 = 1$ as $\omega_i \in \{-1, 1\}$ we have this). Thus

$$\omega_1 = - t_1 \sum_{i=2}^N \omega_i t_i$$

(since $t_1 \in \{-1, 1\}$, $t_1^2 = 1$, we have this). Thus ω_1 is

completely determined by $\{\omega_2, \dots, \omega_N\}$. Thus if we hold $\{\omega_2, \dots, \omega_N\}$ we can't make any change in ω_1 , without violating the constraint ③ in the optimization.

Thus if we want to optimize some subset of the ω_i 's we must update at least two of them simultaneously in order to keep satisfying the constraints. This motivates the SMO algorithm which simply does the

following:

Repeat till convergence {
 1. Select some pair ω_i & ω_j to update next (using a heuristic that tries to pick the two that allow us to make the biggest progress towards global maximum).

2. Re optimize $W(a)$ with respect to a_i, a_j , while holding all other a_k 's ($k \neq i, j$) fixed. (6)
- 3.

To test the convergence of the algorithm we ~~test~~ ^{check} Karush-Kuhn-Tucker (KKT) conditions, whether are satisfied within a tolerance or not. KKT conditions ~~can be~~ are given by

$$a_i = 0 \Rightarrow t_i (w^T x_i + b) \geq 1$$

$$a_i = c \Rightarrow t_i (w^T x_i + b) \leq 1$$

$$0 \leq a_i \leq c \Rightarrow t_i (w^T x_i + b) = 1.$$

See the last page for derivation of these conditions.

The key reason for SMO to be an efficient algorithm is that the update to ~~a_i, a_j~~ ^{a_i, a_j} can be computed very efficiently. The brief sketch of the SMO ~~is as follows~~ for deriving efficient update is as follows.

Suppose we have decided to hold a_3, \dots, a_N fixed and want to re-optimize $W(a_1, \dots, a_N)$ wrt a_1 & a_2 , subject to the constraints (2) & (3).

From (3) we can write

$$a_1 t_1 + a_2 t_2 = - \underbrace{\sum_{i=3}^N a_i t_i}_{\phi}$$

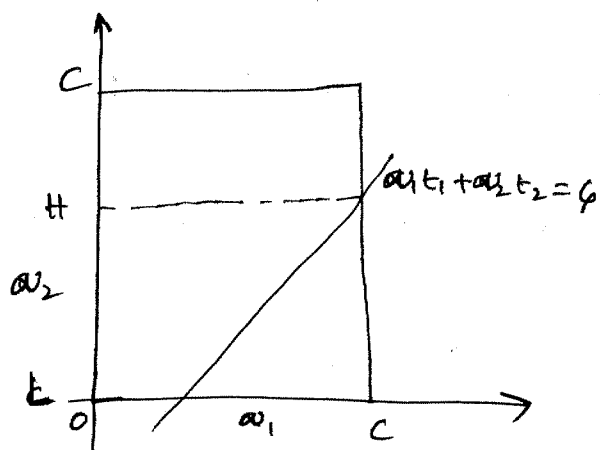
Since RHS is fixed we let denote it by a constant ϕ .

Then

$$a_1 = \frac{\phi - a_2 t_2}{t_1} \quad a_1 t_1 + a_2 t_2 = \phi \quad (4)$$

(7)

We can picture the constraints ω_1 & ω_2 as follows:



~~from~~ ~~from~~

From constraint (2) we know that ~~\omega_1, \omega_2~~ ω_1, ω_2 must lie within the box $[0, C] \times [0, C]$ shown.

also plotted is the line $\omega_1 t_1 + \omega_2 t_2 = \phi$ on which ω_1, ω_2 must lie. We note that from these two constraints $L < \omega_2 < H$.

Now we can rewrite eq (4) as

$$\omega_1 = t_1 (\phi - \omega_2 t_2) \quad (5)$$

Hence the objective $w(\omega)$ can be written as

$$w(\omega_1, \omega_2, \dots, \omega_m) = w(t_1 (\phi - \omega_2 t_2), \omega_2, \omega_3, \dots, \omega_n)$$

Treating $\omega_3, \dots, \omega_m$ as constants, this is just some quadratic function in ω_2 . i.e. this can also be expressed in the form $\alpha \omega_2^2 + \beta \omega_2 + \gamma$ for some appropriate α, β , & γ . If we ignore the box constraint ^(or $L \leq \omega_2 \leq H$) we can

easily maximize this quadratic function by setting its derivative to zero & solving. Let that be $\omega_2^{\text{new, unclipped}}$,

But we want to maximize ' $w(\omega)$ ' ~~not~~ subject to box constraint & $L \leq \omega_2 \leq H$, in that case we have

$$\omega_2^{\text{new}} = \begin{cases} H & \text{if } \omega_2^{\text{new, unclipped}} > H \\ \omega_2^{\text{new, unclipped}} & \text{if } L \leq \omega_2^{\text{new, unclipped}} < H \\ L & \text{if } \omega_2^{\text{new, unclipped}} < L. \end{cases}$$

With this ω_2^{new} now we can use (5) to find ω_1^{new} ; Here few considerations to be made are the choice of the ~~best~~ heuristics used to

Select the next ~~ω_i, ω_j~~ ω_i, ω_j to update and other is $\textcircled{8}$
 how to update 'b' as the SMO algorithm runs.

derivation of Karush-Kuhn-Tucker conditions (KKT):

we have ^{Lagrangian} ~~for~~ the cost function for nonseparable case as

$$L(\omega, b, \xi, \alpha, \gamma) = \frac{1}{2} \omega^T \omega + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [t_i (x_i^T \omega + b) - 1 + \xi_i] - \sum_{i=1}^N \gamma_i \xi_i \quad \text{--- } \textcircled{9}$$

KKT conditions for this are

$$\gamma_i \xi_i = 0 \quad \text{--- } \textcircled{1}$$

$$\alpha_i [t_i (x_i^T \omega + b) - 1 + \xi_i] = 0$$

or

$$t_i (x_i^T \omega + b) - 1 + \xi_i \geq 0 \quad \text{--- } \textcircled{2}$$

and we have

$$\frac{\partial L}{\partial \xi_i} = 0 \quad \Rightarrow \quad C - \alpha_i - \gamma_i = 0$$

$$\text{(i)} \quad \alpha_i = 0 \Rightarrow C - \gamma_i = 0 \quad \gamma_i > 0 \Rightarrow \xi_i = 0$$

$$\Rightarrow t_i (x_i^T \omega + b) - 1 \geq 0$$

$$\Rightarrow t_i (x_i^T \omega + b) \geq 1 \quad \text{--- } \textcircled{1}$$

$$\text{(ii)} \quad 0 < \alpha_i < C \Rightarrow \gamma_i = 0 \Rightarrow \xi_i = 0$$

$$t_i (x_i^T \omega + b) - 1 = 0 \Rightarrow t_i (x_i^T \omega + b) = 1 \quad \text{--- } \textcircled{ii}$$

$$\text{(iii)} \quad \alpha_i = C \Rightarrow \gamma_i = 0 \Rightarrow \xi_i \geq 0$$

$$t_i (x_i^T \omega + b) - 1 + \xi_i = 0$$

$$t_i (x_i^T \omega + b) = 1 - \xi_i \Rightarrow t_i (x_i^T \omega + b) \leq 1 \quad \text{--- } \textcircled{iii}$$

derivation of dual form in non separable case:

The ~~cost~~ Lagrangian cost function to be minimized is given by

$$L(w, b, \xi, \alpha, \gamma) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i - \sum_{i=1}^N \alpha_i [t_i (x_i^T w + b) - 1 + \xi_i] - \sum_{i=1}^N \gamma_i \xi_i \quad - (6)$$

$$\frac{\partial L}{\partial w} = 0 \Rightarrow w - \sum_{i=1}^N \alpha_i t_i x_i = 0 \quad - (1)$$

$$\frac{\partial L}{\partial b} = 0 \Rightarrow \sum_{i=1}^N \alpha_i t_i = 0 \quad - (2)$$

$$\frac{\partial L}{\partial \xi_i} = 0 \Rightarrow C - \alpha_i - \gamma_i = 0 \quad - (3)$$

We have $\alpha_i \geq 0$ & $\gamma_i \geq 0$ from (3) $\gamma_i = C - \alpha_i$

$$\gamma_i \geq 0 \Rightarrow C - \alpha_i \geq 0 \Rightarrow C \geq \alpha_i$$

$$2) \quad \underline{0 \leq \alpha_i \leq C.}$$

KKT conditions are

$$\alpha_i [t_i (x_i^T w + b) - 1 + \xi_i] = 0$$

$$\gamma_i \xi_i = 0$$

$$\xi_i \geq 0 \quad \&$$

$$t_i (x_i^T w + b) - 1 + \xi_i \geq 0$$

Substituting (1), (2), (3) in (6) yields dual form which is to be maximized

$$\text{maximize } W(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j t_i t_j \langle x_i, x_j \rangle$$

$$\text{Subject to } \sum_{i=1}^N t_i \alpha_i = 0 \quad \& \quad 0 \leq \alpha_i \leq C, \quad i=1, \dots, N.$$