**Machine Learning (IS ZC464)** Session 11:

Feed forward Neural Networks – Multilayer Perceptron (MLP) and Radial Basis Function Neural Network (RBFNN)

# Computing Gradient

$E(w)$ = $\sum_i T_i^2$

= $\sum_i (y^i - g(h_w(x^i)))^2$

where $T_i$ is the error term for the i[th] observation and is given by the difference between the desired output ($y^i$) value and the estimated value ($h_w(x^i)$) of the output

$T_i = y^i - g(h_w(x^i))$

$h_w(x^i)$ is the hypothesis function given by

$h_w(x^i) = w_1 x^i_1 + w_2 x^i_2 + \ldots w_n x^i_n$

Observe: E is a function of w.

Note: Here the superscript 'i' represents the 'i'th observation and NOT the power of x.

# Observe

- E is the function of $T_i$  <span>Revision</span>

- Ti is the function of g (assuming y as constant)

- g is the function of h

- h is the function of w

- Chain rule of Differentiation

$$\partial E/\partial w_k = \sum_i \partial E/\partial T_i * \partial T_i/\partial g * \partial g/\partial h * \partial h/\partial w_k$$

<span>Equation 1</span>

# Observe

Since

$$E(w) = \sum_i T_i^2$$

$$\partial E / \partial T_i = 2 * T_i$$

Chain rule of Differentiation

$$\partial E / \partial w_k = 2 * \sum_i T_i * \partial T_i / \partial g * \partial g / \partial h * \partial h / \partial w_k$$

Also since

$$T_i = y^i - g(h_w(x^i))$$

Therefore

$$\partial T_i / \partial g = 0 - 1 = -1$$

Revision

Equation 2

# Working with derivatives

Equation 2 now becomes

$$\partial E/\partial w_k = 2 * \sum_i T_i * (-1) * \partial g/\partial h * \partial h/\partial w_k$$

Also since

$$\partial g/\partial h = \partial g(h_w(x^i))/\partial h = g'$$

Equation 3

And

$$h_w(x^i) = w_1 x^i_1 + w_2 x^i_2 + \ldots w_n x^i_n$$

Therefore

$$\partial h/\partial w_k = x^i_k$$

Hence equation 3 is simplified as

$$\partial E/\partial w_k = -2 * \sum_i T_i * g' * x^i_k$$

Equation 4

# Computing gradient in the direction of $w_k$

- Substitute expression for $T_i$ in equation 4

$$\partial E/\partial w_k = -2 * \sum_i (y^i - g(h_w(x^i))) * g' * x^i_k$$

Equation 5

- The Weight update in the direction of $w_k$

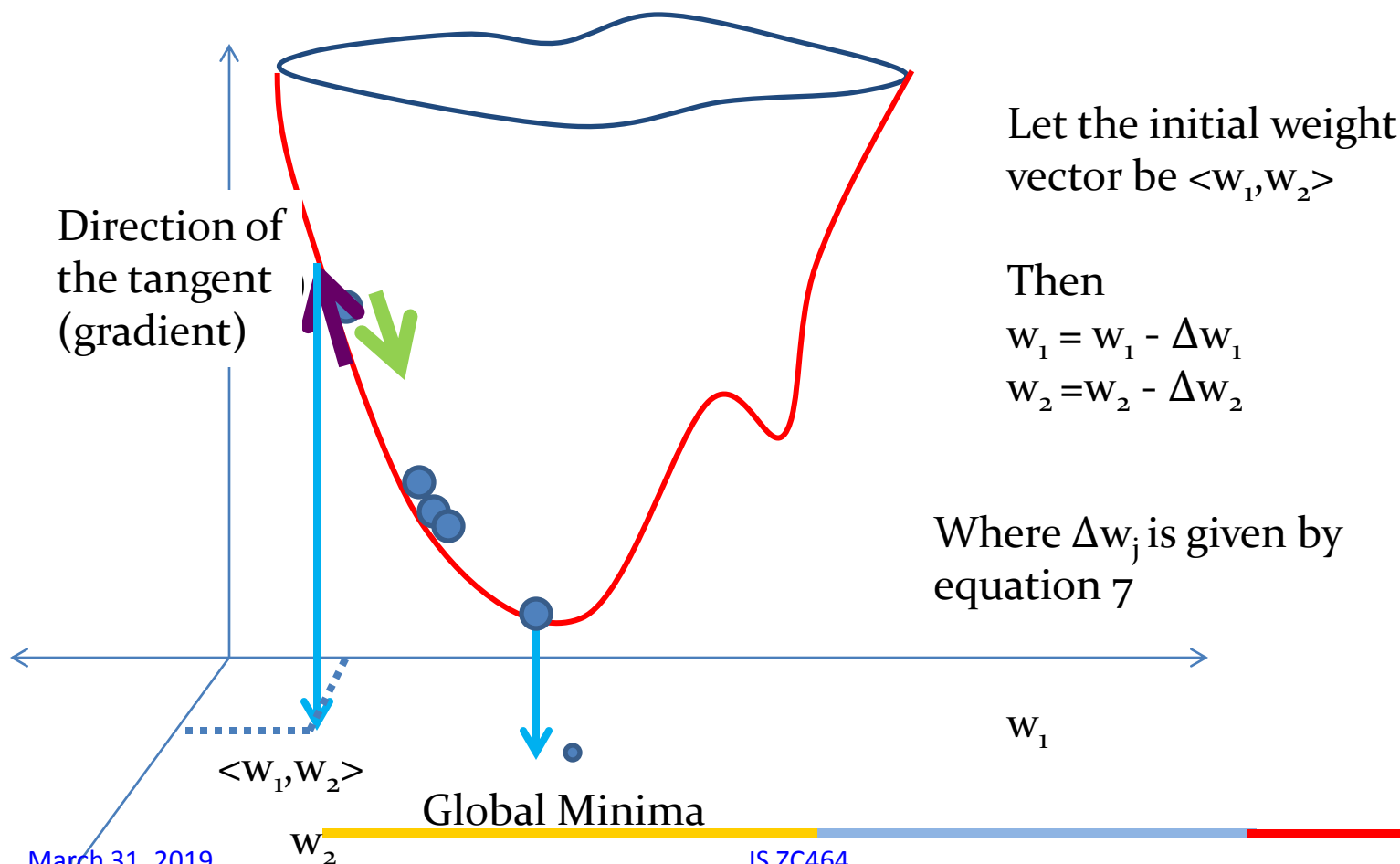$$\Delta w_k = -2 * \sum_i (y^i - g(h_w(x^i))) * g' * x^i_k$$

Equation 6

Where 2 can be dropped to bring normalization.

$$\Delta w_k = - \sum_i (y^i - g(h_w(x^i))) * g' * x^i_k$$

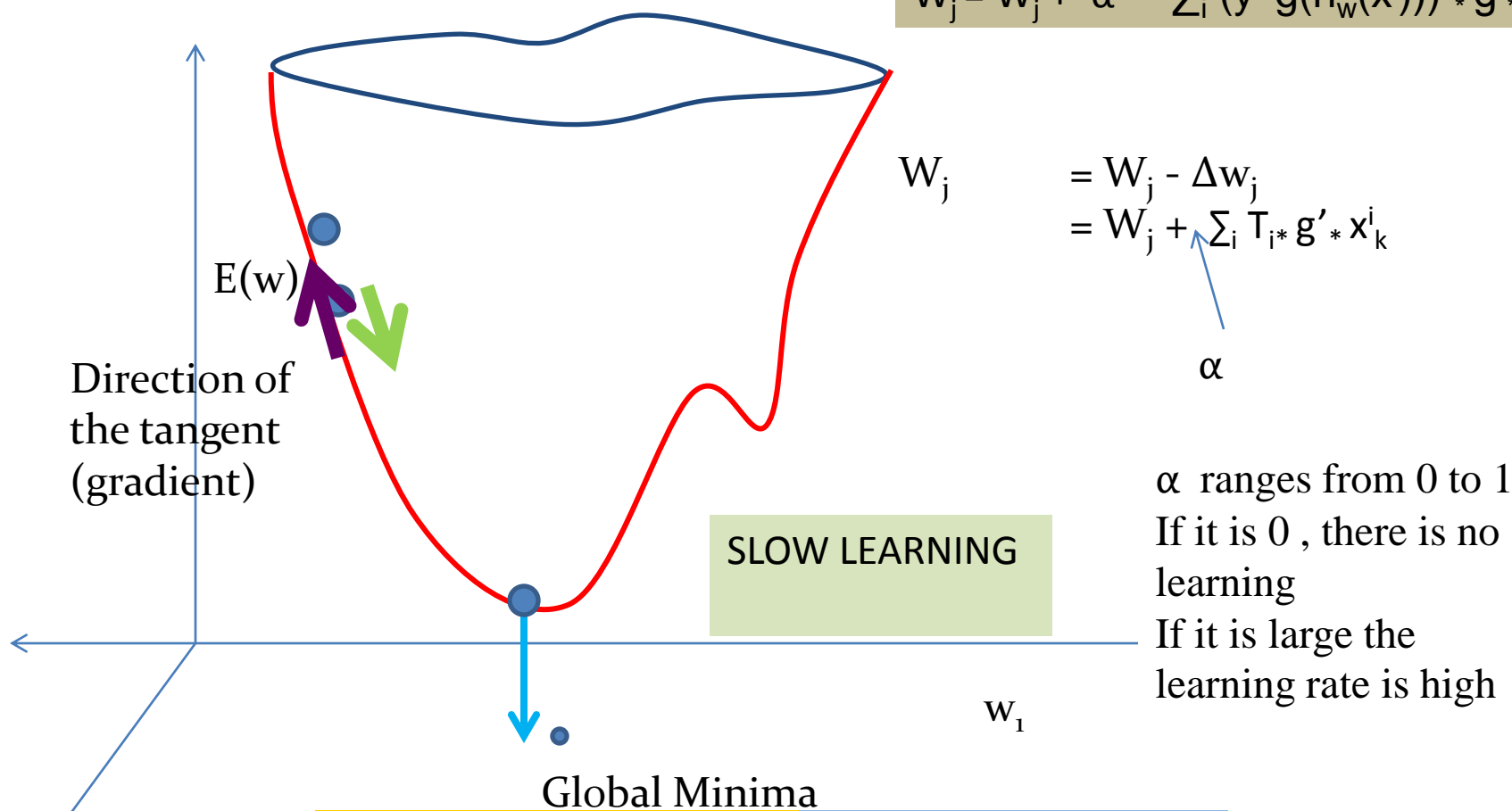Equation 7

# Delta Learning: Modification of the Initial weight

## Revision

Direction of the tangent (gradient)

Let the initial weight vector be $<w_1, w_2>$

Then
$w_1 = w_1 - \Delta w_1$
$w_2 = w_2 - \Delta w_2$

Where $\Delta w_j$ is given by equation 7

$<w_1, w_2>$

$w_1$

Global Minima

$w_2$

# Learning rate: fast or slow learning

$$W_j = W_j + \alpha * \sum_i (y^i - g(h_w(x^i))) * g' * x^i_k$$

$$W_j = W_j - \Delta w_j$$
$$= W_j + \sum_i T_{i*} g' * x^i_k$$

$\alpha$

$E(w)$

Direction of
the tangent
(gradient)

SLOW LEARNING

$\alpha$ ranges from 0 to 1
If it is 0 , there is no
learning
If it is large the
learning rate is high

$w_1$

Global Minima

$w_2$

# Feed Forward and Back propagation

- In feed forward neural network, the weights are computed on the basis of the input propagating through neurons in the forward direction. In this no neuron receives the modified input.

- In back propagation neural network, the processed input is cycled again through the previous layer neurons and the weights are modified.

# Feed Forward Neural Networks

Weights learning is one way

# Back Propagation Neural Networks

Weights learning is cyclic

If Error <= ToleranceLimit
Then
terminate

X1

X2

xn

W1

W2

wn

$\Sigma$

Error E is computed

If Error > ToleranceLimit

Weights are modified

# Activation function should be differentiable for Delta Learning

– $\Delta w_k = - \sum_i (y^i - g(h_w(x^i)))_* g'_* x^i_k$

- The function g' is 0 if g is not differentiable

- Example Activation functions

- Step Function : Not differentiable

- Sigmoid Function : Differentiable

$$y = \frac{1}{1 + \exp(-x)}$$

# Gradient Descent Algorithm

1. Initialize weights in the n-dimensional space randomly.

2. Compute error E.

3. Define error tolerance limit L.

4. While (E > L)

    – Modify weights  W according to delta rule

    – Compute error E with the modified weights and the given input.

# Terminology used in text book

| | Used in the slides here | Used in book by Mitchell (Chapter 4) |
|---|---|---|
| Set of Training samples | Input : vector $x^i$ : i = 1,2,…m <br> output: $y^i$ : i = 1,2,…m | D is the set of training samples $d \in D$ <br> Input : vector $x_d$ : $d \in D$ <br> output : $t_d$ : $d \in D$ |
| Target (Known-supervised) | $y^i$ | $t_d$ |
| Input feature vector | $x^i = <x^i_1, x^i_2, x^i_3, … x^i_n>$ | $x_d = <x_{d1}, x_{d2}, x_{d3}, … x_{dn}>$ |
| Output-predicted by ANN | $h_w(x^i)$ | $o_d$ |
| error | $y^i - h_w(x^i)$ | $t_d - o_d$ |

# What to do with the weights (W) obtained using Gradient Descent Algorithm

- Let W = $<w_1, w_2, w_3, \ldots w_n>$    [as a result of training]

- Have a new feature vector is x = $(x_1, x_2, x_3, \ldots x_n)$ corresponding to the sample not yet seen by the machine (known as test vector)

- Compute output y as follows

- $\quad y = h_w(x) = w_1 x_1 + w_2 x_2 + \ldots w_n x_n$

- This is the identification of the output **[Machine has learned ]**

# Multilayer Feed Forward neural network

- These represent the class of networks which approximate the complex functions.

- The network has one or more hidden layers.

- The neuron 'i' of layer 'L' is connected by a synaptic weight $w_{ki}$ to the 'k'th neuron of layer 'L+1'

# Weight Terminology



$w_{21}$

Layer L

Layer 'L+1'

# MLP

# Multilayer Feed forward Neural Network



$W_{21}$

Hidden Layer

Output Layer

# Multilayer Feed forward Neural Network



$$w_{21}$$

Hidden Layer

Output Layer

# Multilayer Feed forward Neural Network



$w_{21}$

Hidden Layer

Output Layer

# Multi Layer Perceptrons

- These are acyclic directed graphs.

- MLP is a feedforward neural network.

- Can handle non linearly separable data.

- Have different hidden layers of neurons which process the data.

- Training is through weight learning.

- $i^{th}$ layer passes information to i+$1^{th}$ layer

# Real World Problem

- Face Recognition

# A face to recognize
## .......for a Computer

# ... for Humans



Patterns…

Whose face is this ?
Machine can recognize if trained..

Which pattern of numeric values represents a person's face uniquely?

Human's Face Recognition ability is amazing…….

**?**

A

B

C

# A Face Image

- It is simply a grid of numeric values for the machine.

- A machine uses its computational powers to identify patterns from the above numeric values. (Feature Extraction)

- These patterns are unique to a person.

- A face image is represented by various numerical ways such as PCA eigen faces, DCT, wavelets, other statistical methods.

# Understanding Patterns and Pattern Recognition Problem

Training Patterns

Class 1: <1, 2, 3, >       Consecutive integers in ascending order

Class 2: <1,4,9>       Squares of Consecutive integers

Class 3: <-1, -3, -5>       Descending  integers with step size 2

Testing Pattern : <25, 36, 49>

Humans:  Recognize easily (Good Generalization Capability)

Machines: Need Mathematical Models to recognize patterns

# Patterns

- Individual values in the pattern do not give valuable information about the pattern.

- All values in association with each other are informative.

- Patterns have an underlying mathematical structure.

# Complexity of Face Data

- The geometric face features are not robust with respect to variations in expression or illumination conditions.

- Mathematical representations such as coefficients of the Discrete Cosine Transform, Wavelet Transforms etc. are used to represent the face.
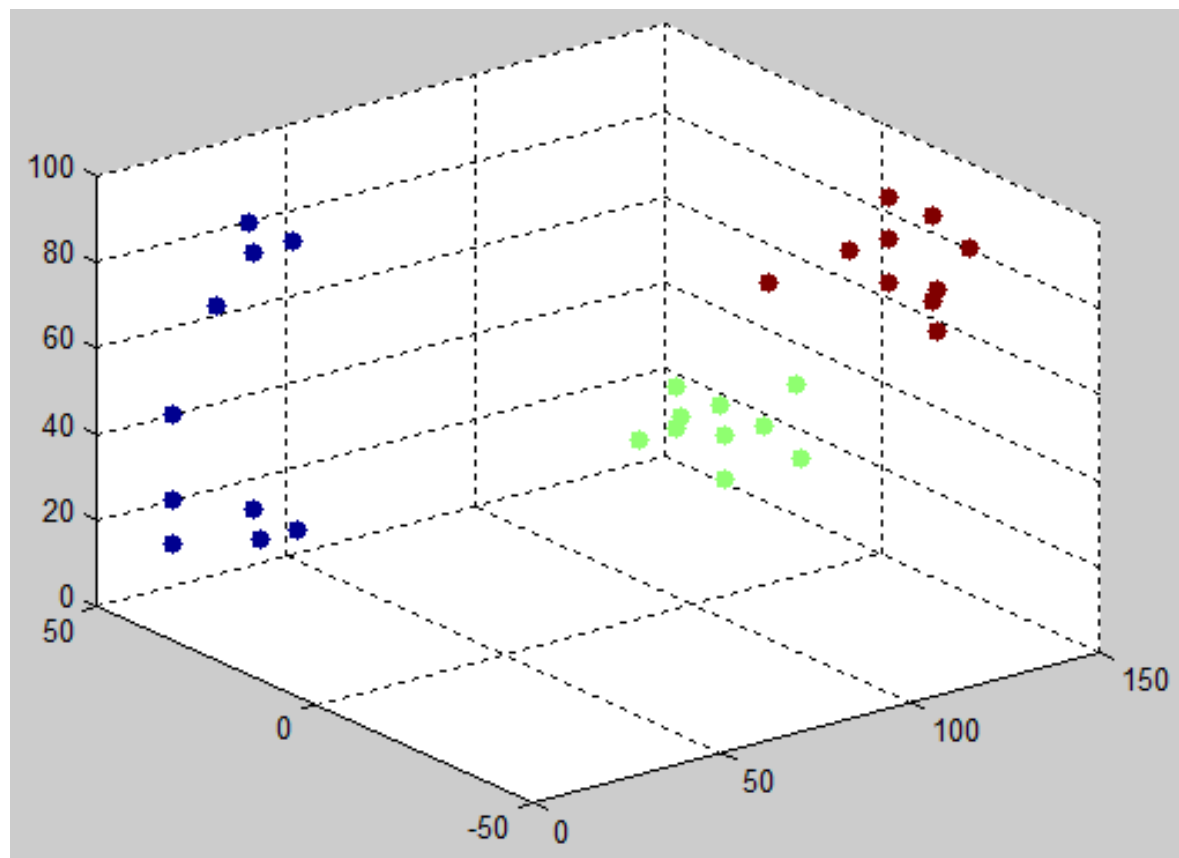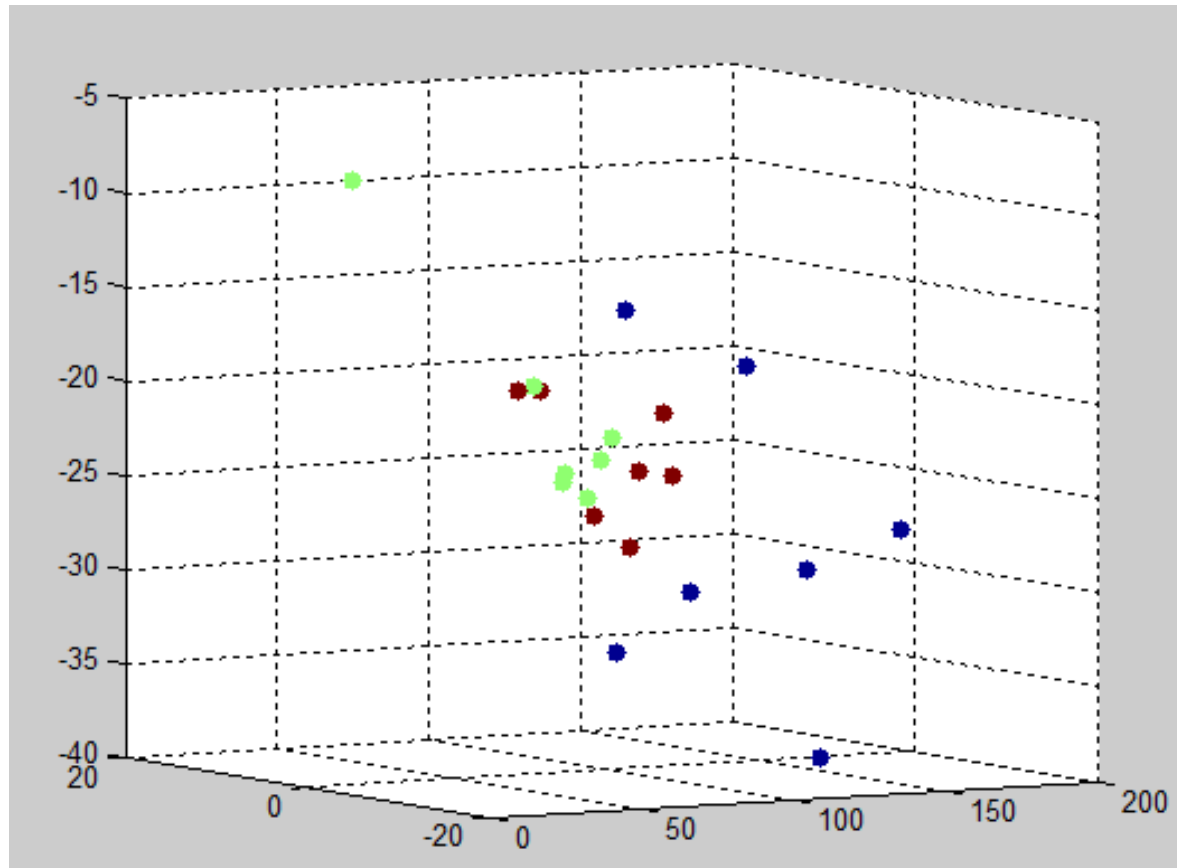
# Complexity of Face Data

- A large number of such coefficients are required to retain identity of a person face.

- A small number of the **Optimal** Features are selected. (to reduce computational load)

- The number (n) of optimal features is also high (e.g. 45 as against all 10000 pixels)
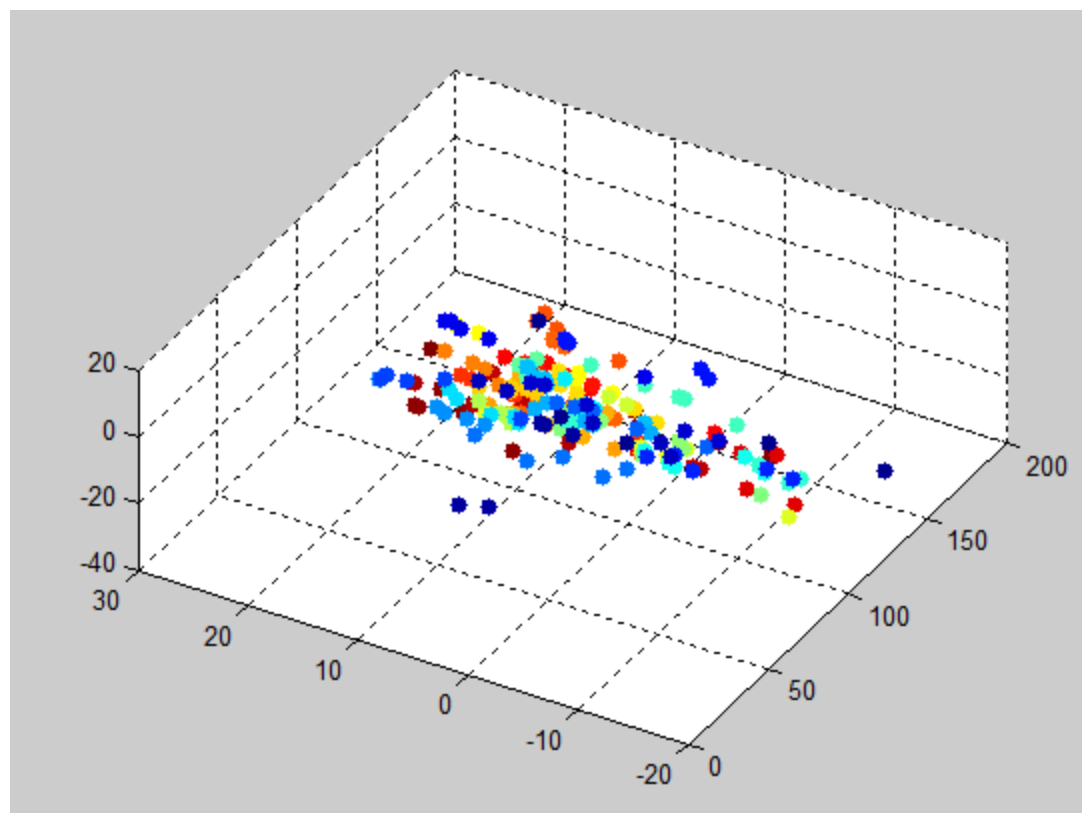
IS ZC464
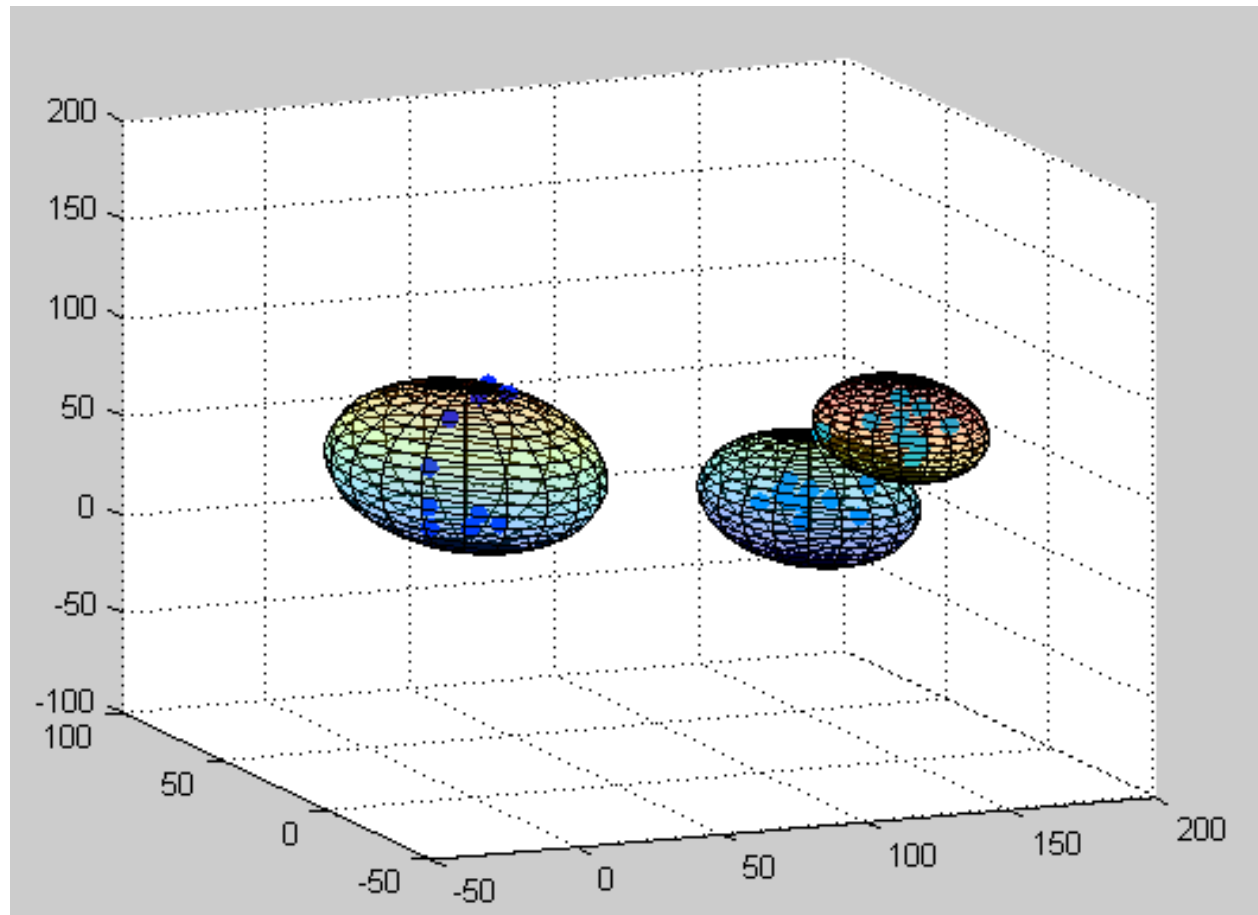
# Linearly Separable Non-Face Data

# Each face is a point in the n-dimensional space. (ORL face data for three persons)

# The points in the n-dimensional space cannot be clustered (colorwise) by hyperplanes.
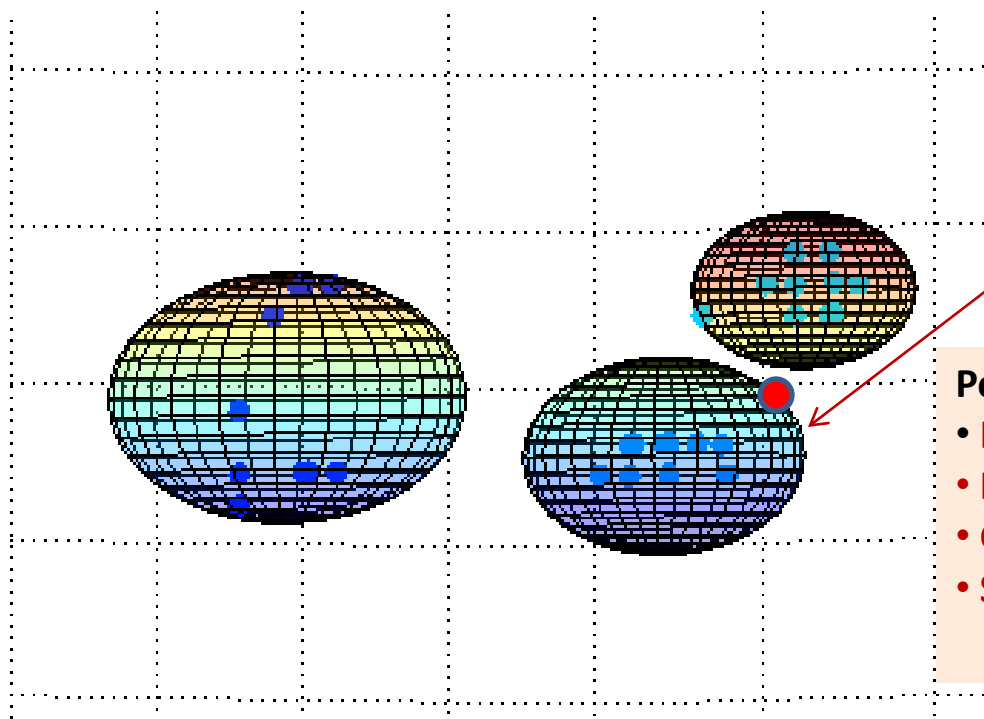
# Face data is nonlinearly separable (Hyper-Surfaces can create boundaries between clusters)

# Classification Problem

Given Training Data



Closest cluster to the n-dimensional test feature vector is computed
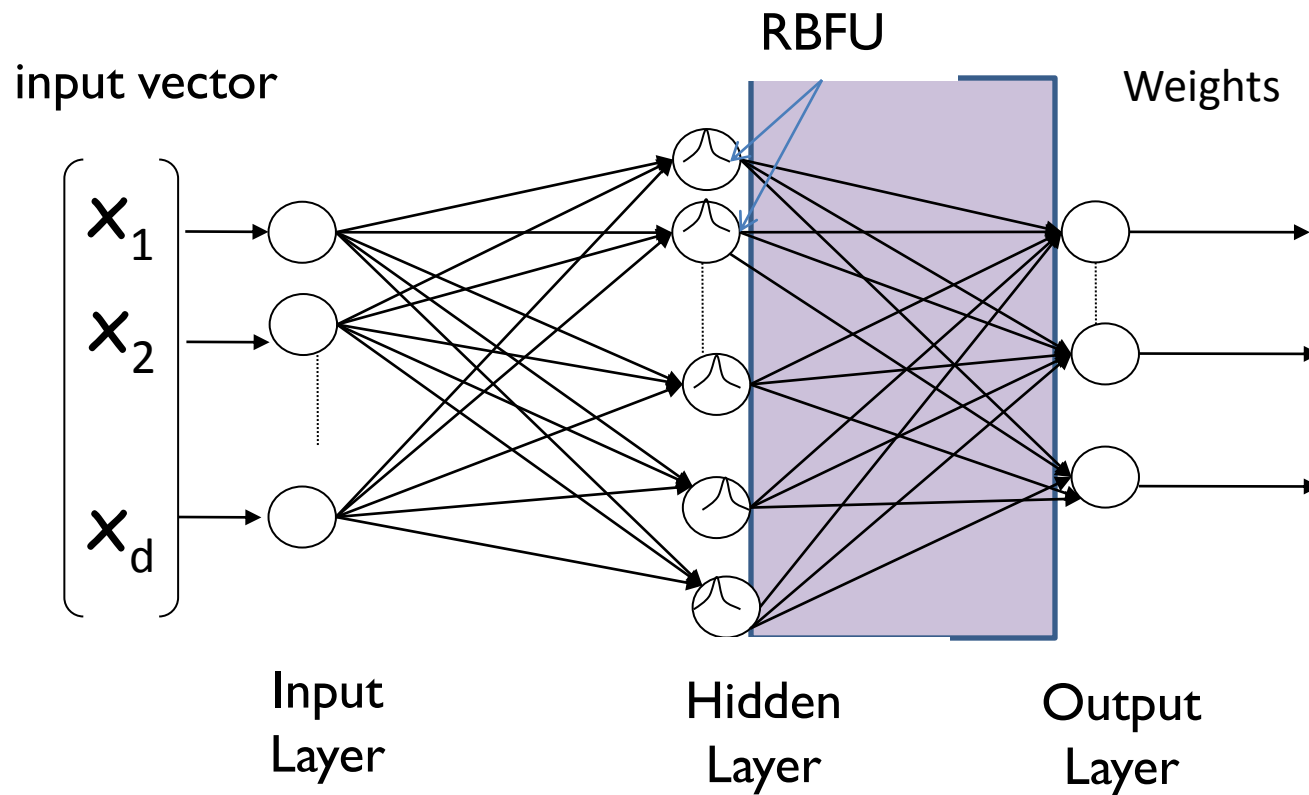
**Possible Decision Boundaries**
- **Hyper Plane**
- **Hyper Sphere**
- **Gaussian Surface**
- **Support Vectors**

**Challenge:** Design of Decision Boundary

# Face Recognition Problem

- Posed as a classification problem

- Classes are the person names (identity)

- Training face images are visualized as points in d-dimensional space (d: pattern size)

- Challenge is in identifying appropriate boundaries demarcating individual cluster.

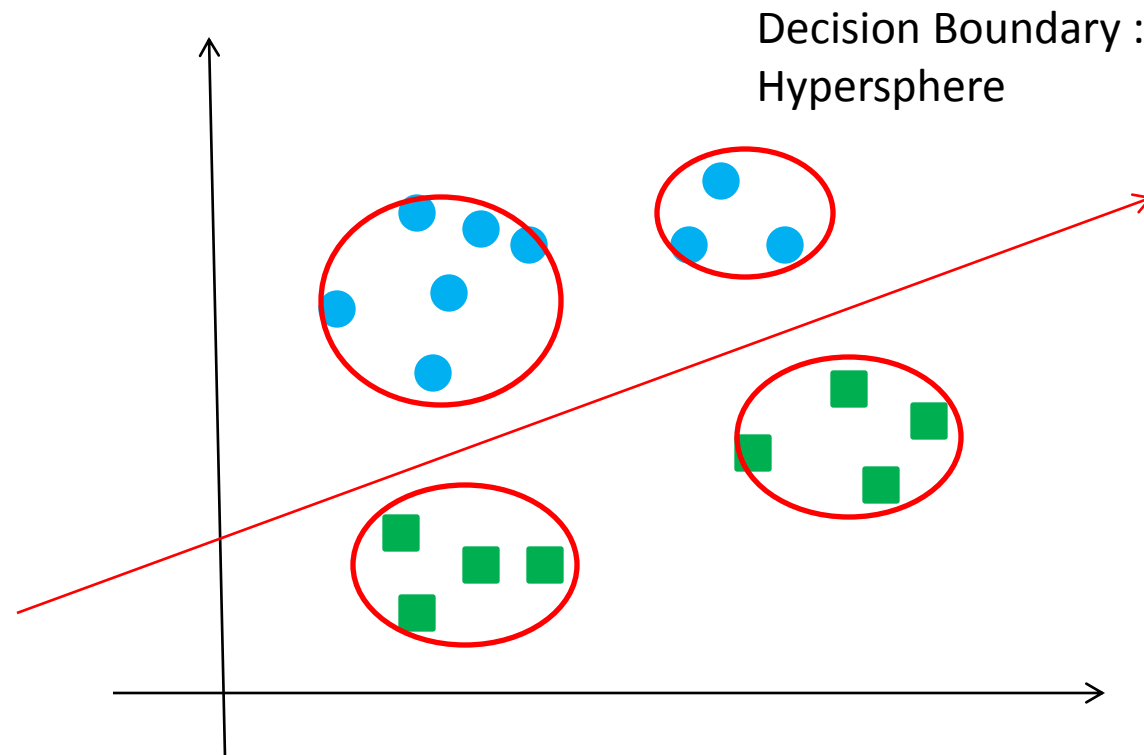# Radial Basis Functions Neural Network (RBFNN) Architecture

# Why More neurons?

- More number of classes.

- Large input sizes of the patterns

- Nonlinear separability of clusters in n-dimensional space.

# RBFNN

- They capture the training environment in terms of weights.

- The radial basis functions units (RBFU) locally capture the structure of the data

- Basis functions at the RBFU play an important role in transforming the nonlinearly separable high dimensional data to a space of linearly separable data.

# Multi Layer Perceptron Vs. RBFNN
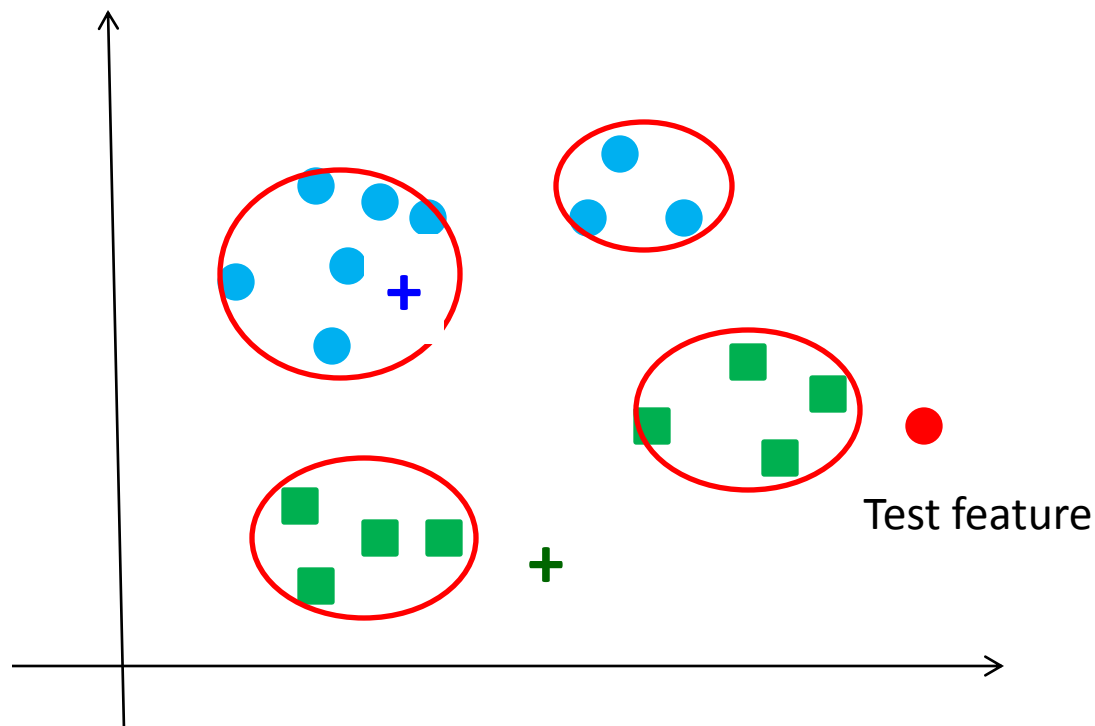


Decision Boundary : Hypersphere

The center of the natural cluster is the center of the hidden neuron

A hidden neuron is sensitive for data points near its center

# Nearest Neighbor Classification Vs.  RBFNN based classification

"Do not know " condition can be handled well by RBFNN

Test feature

**Nearest neighbor:**  Shortest distance to the mean of the cluster

**RBFNN:** Within limits of Radial distance to the mean of the cluster