# Support Vector Machines
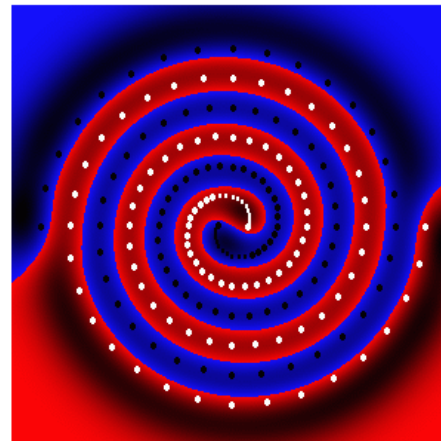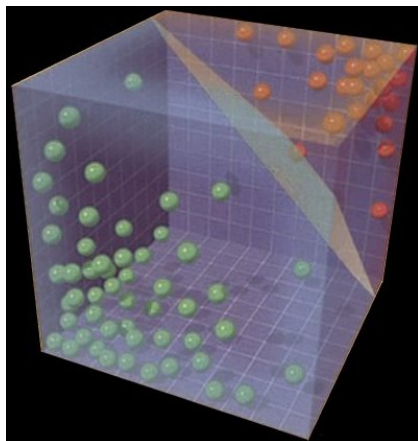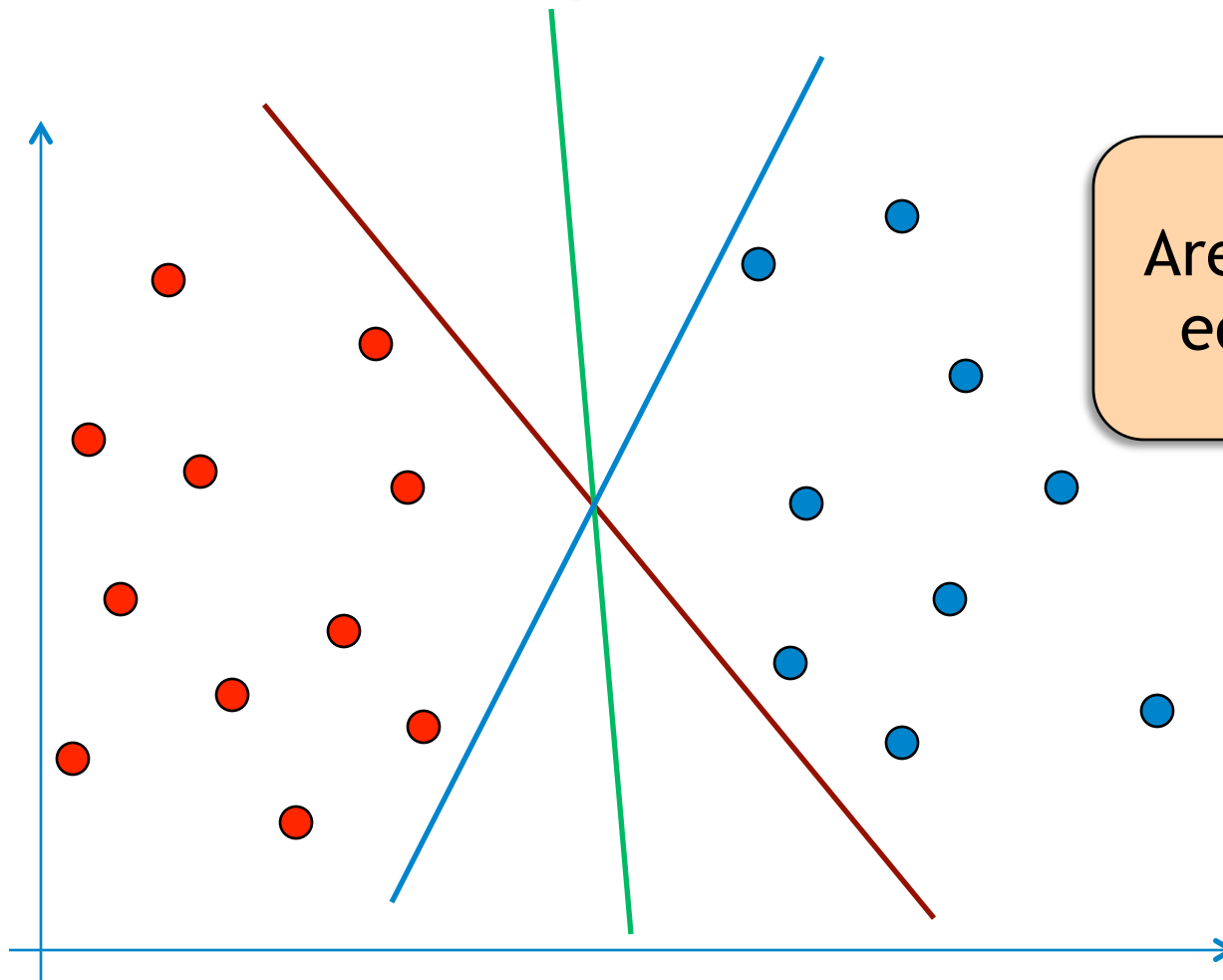## Maximizing the Margin

Anoop M. Namboodiri

Centre for Visual Information Technology

IIIT, Hyderabad, INDIA

IIIT Hyderabad

# Perceptron Learning



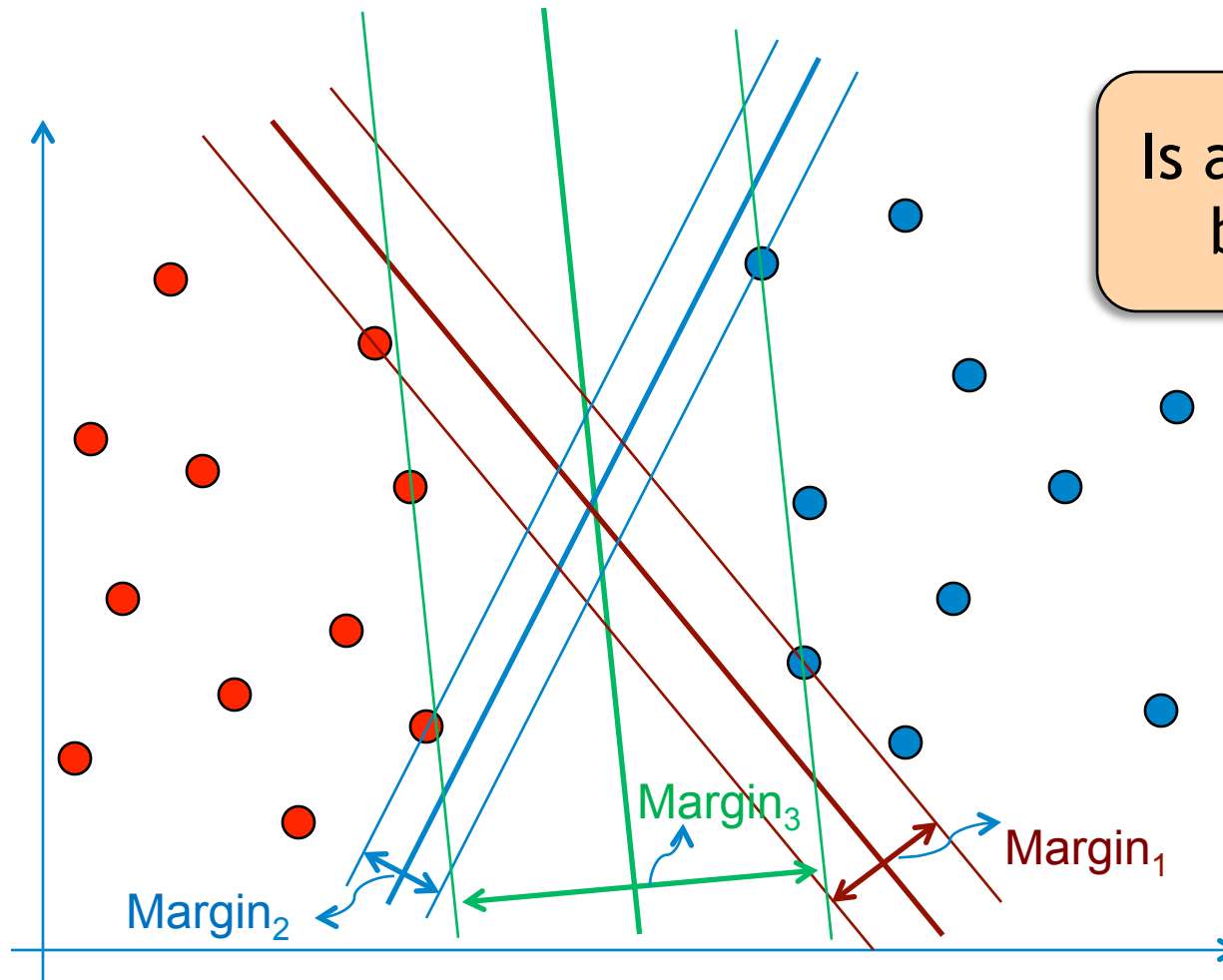Are all solutions equally good?

- Multiple solutions exist for linearly separable data
- Perceptron learning (any GD) results in a feasible solution

# Margin: The No-mans Band



Is a Larger Margin better? Why?

Margin$_3$

Margin$_1$

Margin$_2$

- Margin: Width of a band around decision boundary without any training samples
- Margin varies with the position and orientation of the separating hyperplane

IIIT Hyderabad

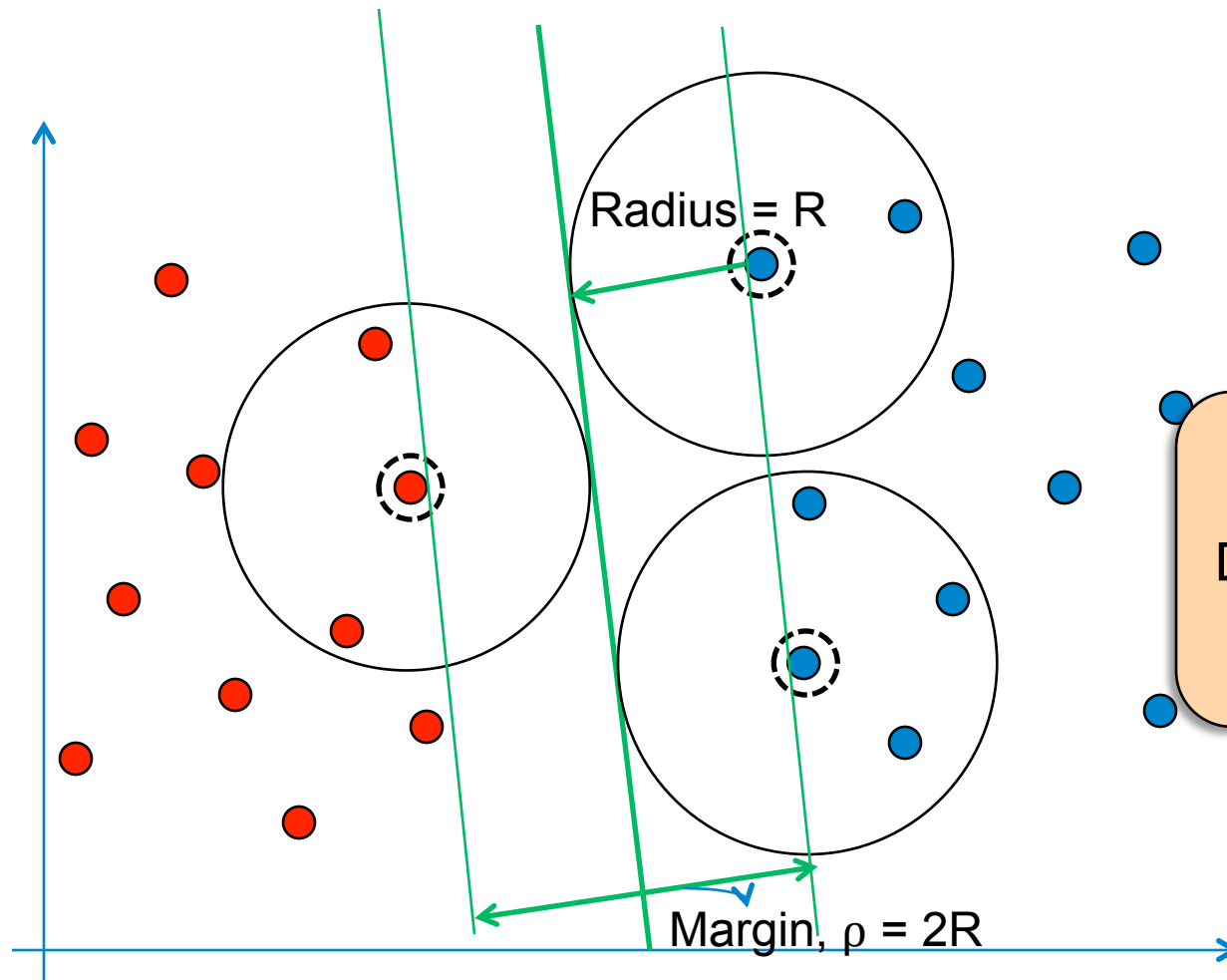A few samples control the Dec. boundary

- Margin: Radius of a region around each training sample, through which the decision boundary cannot pass

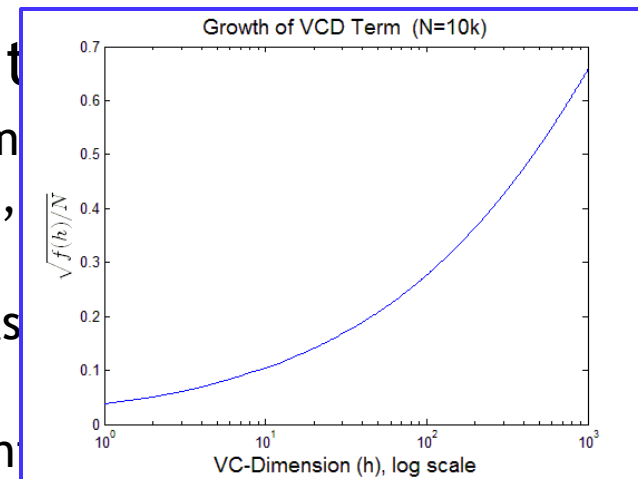- As margin increases, the feasible region reduces

# Margin: Band vs. Bubbles



Radius = R

Samples that support the Decision boundary are called *Support Vectors*

Margin, $\rho = 2R$

- Margin: Both interpretations yield the same decision boundary

IIIT Hyderabad

# Can we Quantify Generalization?

- Work by Vapnik and Chervonenkis in t

    1. V.N. Vapnik, A.Ya. Červonenkis, "On Uniform
       Frequencies of Events to their Probabilities",
       Primenen, 1971, 16(2), pp. 264-279.
    2. V.N. Vapnik, "Estimation of Dependences Bas
       Russian]. Nauka, Moscow, 1979.
    3. V.N. Vapnik, "The Nature of Statistical Learn
       Verlag, New York, 1995.



Growth of VCD Term (N=10k)

- Bound on Expected loss [3]: $R(\alpha) \le R_{train}(\alpha) + \sqrt{f(h)/N}$

- $h$ is the VC dimension, and $f(h)$ is given by:

$$f(h) = h + h\log(2N) - h\log(h) - c$$

# Why Maximize the Margin?

- To reduce test error, keep training error low (say 0), and minimize the VC-dimension, $h$.

$$\text{Relative Margin}: \frac{\rho}{D}$$

$$\text{VC - D}, \ h \le \min\left\{ d \ , \ \left\lceil \frac{D^2}{\rho^2} \right\rceil \right\} + 1$$

- Maximizing margin improves generalization.
- $h$ can be made independent of the dimensionality: $d$.

Data dia: D

Margin: ρ

# Formalizing the Margin

Dec. Boundary: $W^TX + b = 0$

Let parallel hyperplanes be:
$W^TX + b = \pm\varepsilon$

- Note: The value of $W^TX_i + b$ is dependent on the scale of X and W

Anoop M. Namboodiri

8

# Formulation

- Let $g(x) = \mathbf{w}^T\mathbf{x} + b$.
- We want to maximize k such that:
  - $\mathbf{w}^T\mathbf{x}_i + b \geq k$  for  $d_i = 1$
  - $\mathbf{w}^T\mathbf{x}_i + b \leq -k$  for  $d_i = -1$

- Value of $g(\mathbf{x})$ dependents on $\|\mathbf{w}\|$ :
  1. Keep $\|\mathbf{w}\| = 1$, and maximize $g(\mathbf{x})$, or
  2. Let $g(\mathbf{x}) \geq 1$, and minimize $\|\mathbf{w}\|$.

- We use approach (2) and formulate the problem as:
  - Minmize:    $\frac{1}{2}\mathbf{w}^T\mathbf{w}$
  - Subject to: $d_i(\mathbf{w}^T\mathbf{x}_i + b) \geq 1$,  for $i = 1..N$



$g(\mathbf{x}) = 0$

$g(\mathbf{x}) = -k$

$g(\mathbf{x}) = k$

$\mathbf{w}$

# The Optimization Problem

Minimize : $\Phi(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{w}$

Subject to : $d_i(\mathbf{w}^T\mathbf{x_i} + b) - 1 \geq 0 \quad \forall i$

Saddle Point



- Quadratic objective function with linear inequalities as constraints: QP Solver.

- Integrating the constraints into the Lagrangian form, we get:

Minimize : $J(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2}\mathbf{w}^T\mathbf{w} - \sum_{i=1}^{N}\alpha_i d_i(\mathbf{w}^T\mathbf{x_i} + b) + \sum_{i=1}^{N}\alpha_i$

Subject to : $\alpha_i \geq 0 \quad \forall i$

- Minimize $J$ with respect to $w$ and $b$, and maximize with respect to $\alpha$.

# Converting to the Dual Form

Objective: $J(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^{N} \alpha_i d_i (\mathbf{w}^T \mathbf{x_i} + b) + \sum_{i=1}^{N} \alpha_i$

At the optimum, $1: \dfrac{\partial J}{\partial \mathbf{w}} = 0$ and $2: \dfrac{\partial J}{\partial b} = 0$

KKT Conditions

$1: \mathbf{w}_o = \sum_{i=1}^{N} \alpha_i d_i \mathbf{x_i}$ 

$2: \sum_{i=1}^{N} \alpha_i d_i = 0$ 

$3: \alpha_i [d_i (\mathbf{w}_o^T \mathbf{x_i} + b_o) - 1] = 0$

Obj: $J(\mathbf{w}, b, \boldsymbol{\alpha}) = \sum_{i=1}^{N} \alpha_i + \frac{1}{2} \mathbf{w}^T \mathbf{w} - \mathbf{w}^T \sum_{i=1}^{N} \alpha_i d_i \mathbf{x_i} - b \sum_{i=1}^{N} \alpha_i d_i$

Using 1,2: $Q(\boldsymbol{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x_i}^T \mathbf{x}_j$

# Solving the Dual Form

$$Q(\boldsymbol{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x_i}^T \mathbf{x}_j$$

$$\text{Subject to } \alpha_i \geq 0 \quad \forall_i \quad \text{and} \quad \sum_{i=1}^{N} \alpha_i d_i = 0$$

QP Solver

$\alpha_i$

- The only unknowns (variables) are $\alpha_i$s.
- The constraints are also on $\alpha_i$s only.
- Data vectors appear only as dot products
- Objective is convex, subject to linear constraints
- Can be solved using standard convex quadratic program solvers

$$\mathbf{w}_o = \sum_{i=1}^{N} \alpha_i d_i \mathbf{x_i}$$

KKT Conditions

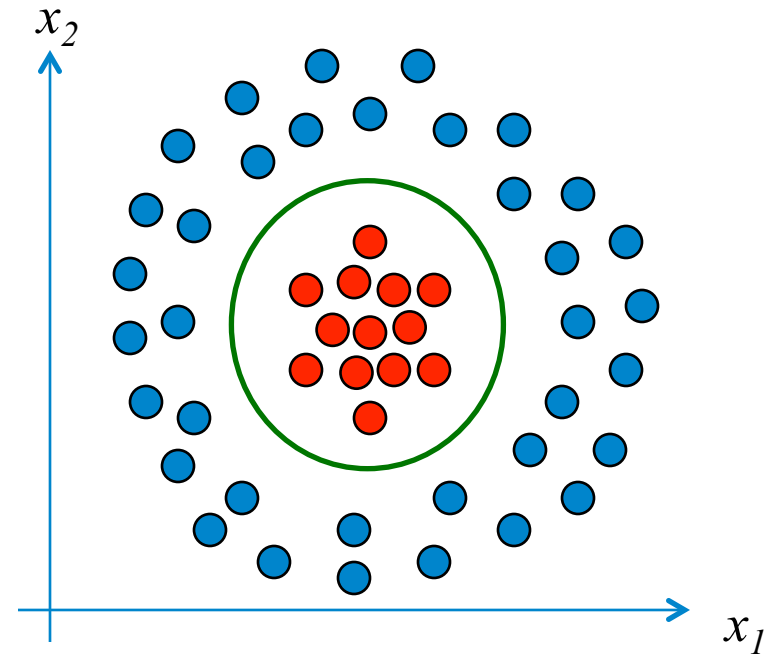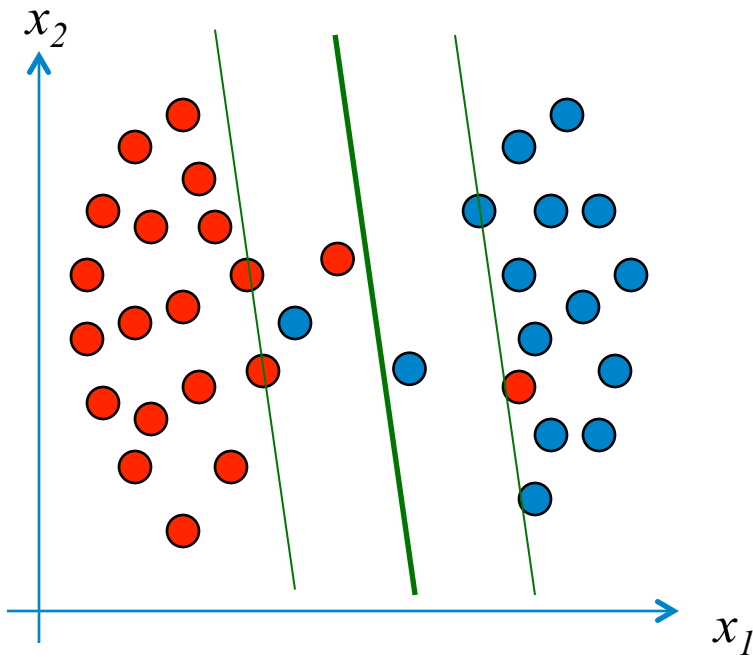$$\alpha_i [d_i (\mathbf{w}_o^T \mathbf{x_i} + b_o) - 1] = 0$$

$$b_o = 1 - \mathbf{w}_o^T \mathbf{x}_{s+}$$

IIIT Hyderabad

Anoop M. Namboodiri

# Non-Separable Data

**1: Noisy Data/Bad Features**

**2: Non-linear Boundary**

Minimize : $\Phi(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{w}$

Subject to : $d_i(\mathbf{w}^T\mathbf{x_i} + b) \geq 1 \quad \forall i$

Introduce slack variables $\quad \xi_i \geq 0$

Minimize : $\Phi(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{w}$

Subject to : $d_i(\mathbf{w}^T\mathbf{x_i} + b) \geq 1 - \xi_i \quad \forall i$



$\xi_i = 0.2$

$\xi_i = 1.7$

$\xi_i = 0.5$

$\xi_i = 2.3$

Also minimize training error $\quad \sum_{i=1}^{N} \mathrm{I}(\xi_i \geq 1) \quad$ or $\quad \sum_{i=1}^{N} \xi_i$

Minimize : $\Phi(\mathbf{w}) = \frac{1}{2}\mathbf{w}^T\mathbf{w} + C\sum_{i=1}^{N}\xi_i$

Subject to : $d_i(\mathbf{w}^T\mathbf{x_i} + b) \geq 1 - \xi_i ; \quad \xi_i \geq 0 , \quad \forall i$

# Dual form with Slack

- Forming the Lagrangian and converting to dual, we get:

$$Q(\pmb{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \mathbf{x_i}^T \mathbf{x}_j$$

$$\text{Subject to } 0 \leq \alpha_i \leq C \quad \forall_i \quad \text{and} \quad \sum_{i=1}^{N} \alpha_i d_i = 0$$

- Note that neither the slack variables, nor their Lagrange multipliers appear in the dual.
- The only change is the additional constraint on $\alpha_i$
- The parameter C controls the relative weight between training error and the VC dimension.

IIIT Hyderabad

# Non-Linear Boundaries

$$x_3 = x_1{}^2 + x_2{}^2$$

$\Phi(x)$

$\Phi$ is a non-linear mapping into a possibly high-dimensional space

- Data vectors occur only as dot products in SVM-learning and testing

Training Phase

- 

$$\text{Label} = sign(\mathbf{w_o} \bullet \Phi(\mathbf{x_{test}}) + b_o)$$

$$\mathbf{w}_o = \sum_{i=1}^{N} \alpha_i d_i \, \Phi(\mathbf{x_i})$$

$$\therefore \text{L} \quad Q(\mathbf{\alpha}) = \sum_{i=1}^{N} \alpha_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \alpha_i \alpha_j d_i d_j \, \mathrm{K}(\mathbf{x_i}, \mathbf{x}_j)$$

$$\text{Label} = sign\left( \sum_{i=1}^{N} (\alpha_i d_i \, \mathrm{K}(\mathbf{x_i}, \mathbf{x_{test}})) + b_o \right)$$

# A Simple Quadratic Kernel

$$\text{Let } \Phi(\mathbf{X}) = \Phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1 x_1 \\ x_1 x_2 \\ x_2 x_1 \\ x_2 x_2 \end{bmatrix}$$

We can compute K(X,Y)=(X.Y)² instead of mapping with **Φ** explicitly and then computing dot product.

$$\text{Let } K(\mathbf{X}, \mathbf{Y}) = \Phi(\mathbf{X}) \bullet \Phi(\mathbf{Y}) = \begin{bmatrix} x_1^2 \\ x_1 x_2 \\ x_2 x_1 \\ x_2^2 \end{bmatrix} \bullet \begin{bmatrix} y_1^2 \\ y_1 y_2 \\ y_2 y_1 \\ y_2^2 \end{bmatrix}$$

$$= x_1^2 y_1^2 + 2 x_1 x_2 y_1 y_2 + x_2^2 y_2^2 = \left(x_1 y_1 + x_2 y_2\right)^2 = \left(\mathbf{X} \bullet \mathbf{Y}\right)^2$$

# Similarly for a Cubic Kernel

- Original Space: 2-dimensional

$$\text{Let } K(\mathbf{X}, \mathbf{Y}) = (\mathbf{X} \bullet \mathbf{Y})^3 = (x_1 y_1 + x_2 y_2)^3$$

$$\Phi(\mathbf{X}) = \Phi\left(\begin{bmatrix} x_1 \\ x_2 \end{bmatrix}\right) = \begin{bmatrix} x_1^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_2^3 \end{bmatrix}$$

- Equivalent to working in an 4-dimensional space

- 4:× and 1:+, instead of 16:× and 3:+

Anoop M. Namboodiri

# Similarly for a Cubic Kernel

- Original Space: 3-dimensional

$$\text{Let } K(\mathbf{X}, \mathbf{Y}) = (\mathbf{X} \bullet \mathbf{Y})^3 = (x_1 y_1 + x_2 y_2 + x_3 y_3)^3$$

$$\Phi(\mathbf{X}) = \Phi\left(\begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}\right) = \begin{bmatrix} x_1^3 \\ x_2^3 \\ x_3^3 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_1^2 x_3 \\ x_1 x_3^2 \\ x_2^2 x_3 \\ x_2 x_3^2 \\ x_1 x_2 x_3 \end{bmatrix}$$

- Equivalent to working in an 10-dimensional space
- 5:× and 2:+, instead of 38:× and 9:+

# A Generic Polynomial Kernel

- Adding two Kernels gives you a new Kernel:

$$K(\boldsymbol{X}, \boldsymbol{Y}) = K_1(\boldsymbol{X}, \boldsymbol{Y}) + K_2(\boldsymbol{X}, \boldsymbol{Y}) \quad : \quad \Phi(\boldsymbol{X}) = \begin{bmatrix} \Phi_1(\boldsymbol{X}) \\ \Phi_2(\boldsymbol{X}) \end{bmatrix}$$

$$K(\boldsymbol{X}, \boldsymbol{Y}) = \Phi(\boldsymbol{X}) \bullet \Phi(\boldsymbol{Y}) = \Phi_1(\boldsymbol{X}) \bullet \Phi_1(\boldsymbol{Y}) + \Phi_2(\boldsymbol{X}) \bullet \Phi_2(\boldsymbol{Y})$$

$$K_p(\boldsymbol{X}, \boldsymbol{Y}) = (1 + \boldsymbol{X} \bullet \boldsymbol{Y})^p = 1 + \boldsymbol{X} \bullet \boldsymbol{Y} + (\boldsymbol{X} \bullet \boldsymbol{Y})^2 + \cdots + (\boldsymbol{X} \bullet \boldsymbol{Y})^p$$

- Just adding a 1 and raising to the power of p maps the input vector into a space containing all original dimensions, all 2-products, 3-products,…,p-products.

# Mercer's Theorem

- Using Kernels, we avoid explicit mapping with Φ
- In fact, we do not even have to know what Φ is as long as we are sure there exists a valid Φ.
- Mercer's Theorem:

Any given kernel can be expanded as a series :

$$K(\mathbf{X}, \mathbf{Y}) = \sum_{i=1}^{\infty} \lambda_i \, \Phi_i(\mathbf{X}) \bullet \Phi_i(\mathbf{Y}), \quad \lambda_i > 0; \textit{iff}$$

$K(\mathbf{X}, \mathbf{Y})$ satisfies the Mercer's conditions (symmetric, continuous, positive semi - definite)

# Popular Kernels

- Polynomial:

$$K_p(\mathbf{X}, \mathbf{Y}) = (1 + \mathbf{X} \bullet \mathbf{Y})^p$$

- Radial Basis Function (RBF) or Gaussian:

$$K_r(\mathbf{X}, \mathbf{Y}) = e^{-\frac{1}{2\sigma^2}\|\mathbf{X}-\mathbf{Y}\|_2^2}$$

- Hyperbolic Tangent:

$$K_s(\mathbf{X}, \mathbf{Y}) = \tanh(\beta_0 \mathbf{X} \bullet \mathbf{Y} + \beta_1)$$

IIIT Hyderabad