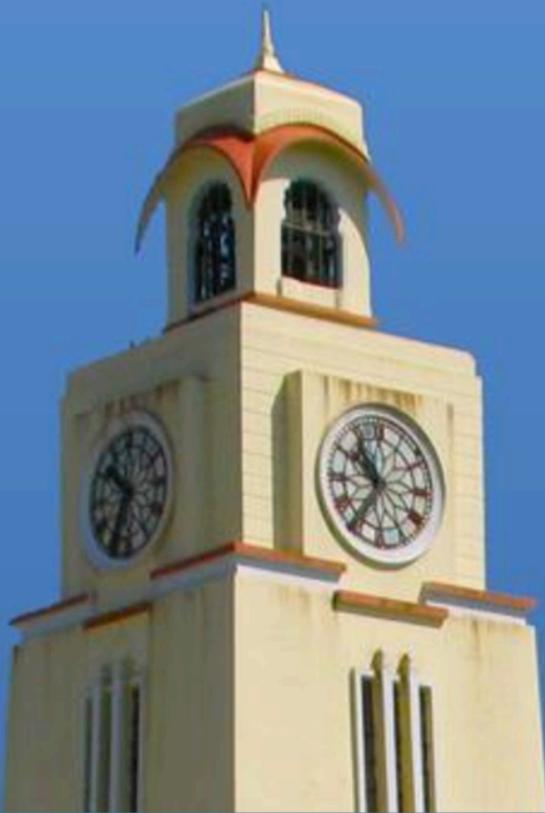




BITS Pilani
Pilani Campus

Support Vector Machines

MFML Team



Text Book(s)

T1	Christopher Bishop: Pattern Recognition and Machine Learning, Springer International Edition
T2	Tom M. Mitchell: Machine Learning, The McGraw-Hill Companies, Inc..

These slides are prepared by the instructor, with grateful acknowledgement of Prof. Tom Mitchell, Prof. Burges, Prof. Andrew Moore and many others who made their course materials freely available online.



Topics to be covered

- Nonlinear SVM
 - Kernel Trick
 - SVM Kernels
 - XOR problem using non linear SVM
-

Optimization Problem

Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} = \frac{1}{2} \|\mathbf{w}\|^2$ is minimized;

and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

$$L(w, b, \alpha_i) = \frac{1}{2} \|w\|^2 - \sum \alpha_i [y_i (w^T x_i + b) - 1]$$

Find $\alpha_1 \dots \alpha_N$ such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \left(\sum_i \sum_j \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j \right)$ is

maximized and

$$(1) \sum \alpha_i y_i = 0$$

$$(2) \alpha_i \geq 0 \text{ for all } \alpha_i$$

Support Vectors

Using KKT conditions :

$$\alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1] = 0$$

For this condition to be satisfied
either $\alpha_i = 0$ and $y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 > 0$

OR

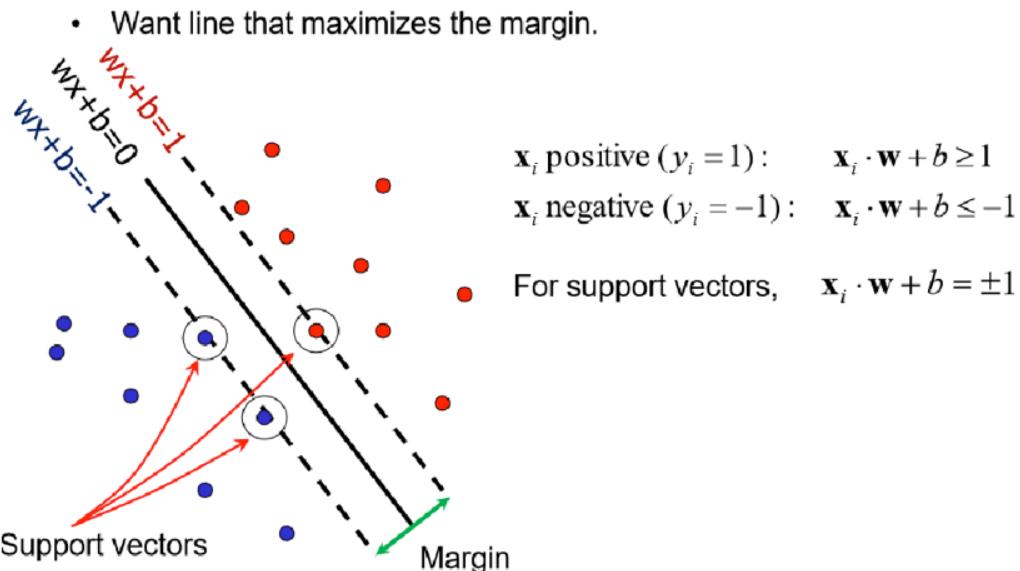
$$y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 = 0 \text{ and } \alpha_i > 0$$

For support vectors:

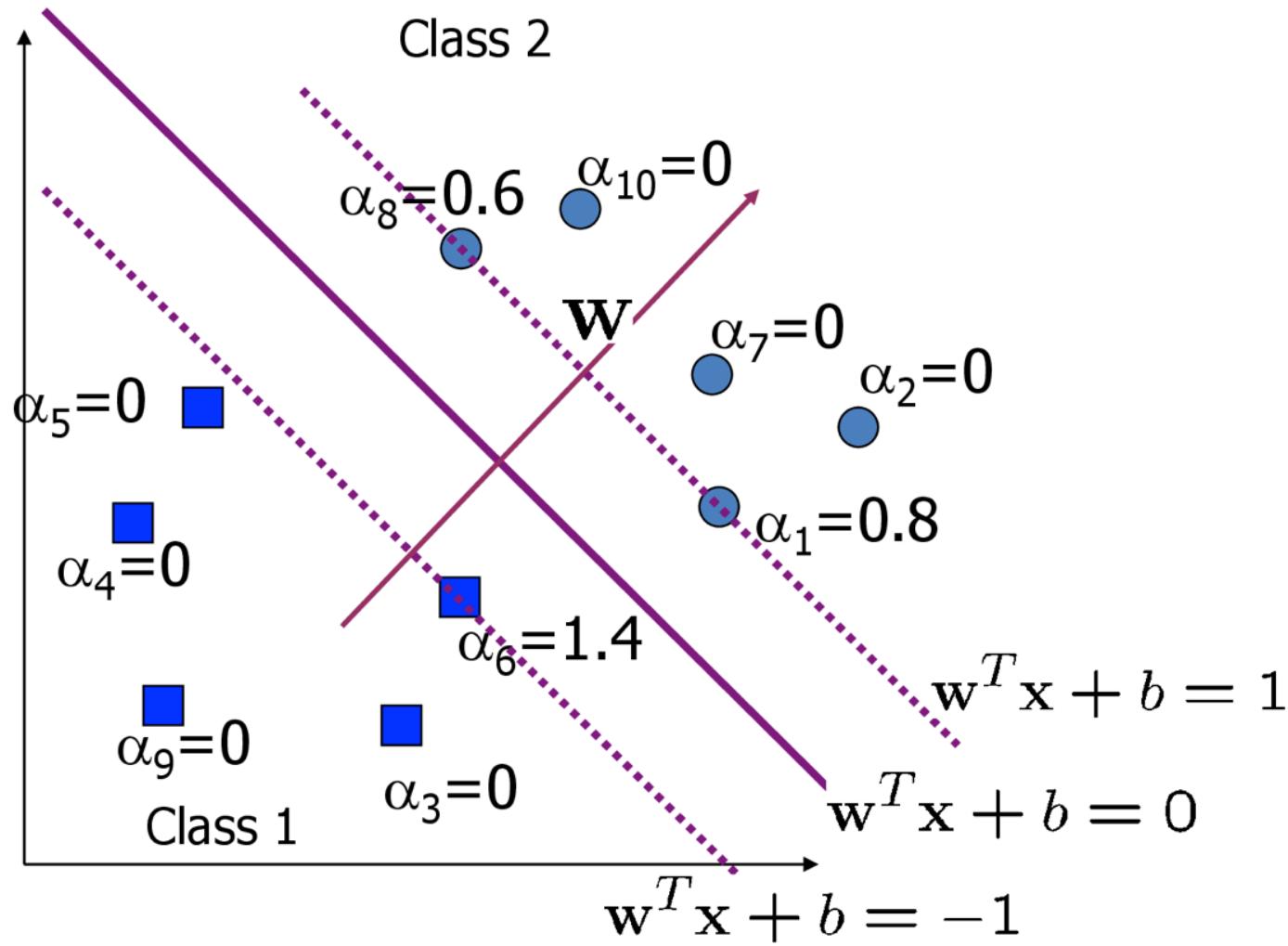
$$y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 = 0$$

For all points other than
support vectors:

$$\alpha_i = 0$$



A Geometrical Interpretation



Solving the Optimization Problem

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$



Solving the Optimization Problem

- Solution: $\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$
 $b = y_i - \mathbf{w} \cdot \mathbf{x}_i$ (for any support vector)

- Classification function:

$$\begin{aligned}f(x) &= \text{sign}(\mathbf{w} \cdot \mathbf{x} + b) \\&= \text{sign}\left(\sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b\right)\end{aligned}$$

If $f(x) < 0$, classify as negative, otherwise classify as positive.

- Notice that it relies on an *inner product* between the test point \mathbf{x} and the support vectors \mathbf{x}_i
- (Solving the optimization problem also involves computing the inner products $\mathbf{x}_i \cdot \mathbf{x}_j$ between all pairs of training points)

Linear SVMs: Overview

- The classifier is a *separating hyperplane*.
- Most “important” training points are support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points x_i are support vectors with non-zero Lagrangian multipliers α_i .
- Both in the dual formulation of the problem and in the solution training points appear only inside dot products:

Find $\alpha_1 \dots \alpha_N$ such that

$Q(\alpha) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j x_i^T x_j$ is maximized and

$$(1) \sum \alpha_i y_i = 0$$

$$(2) 0 \leq \alpha_i \leq C \text{ for all } \alpha_i$$

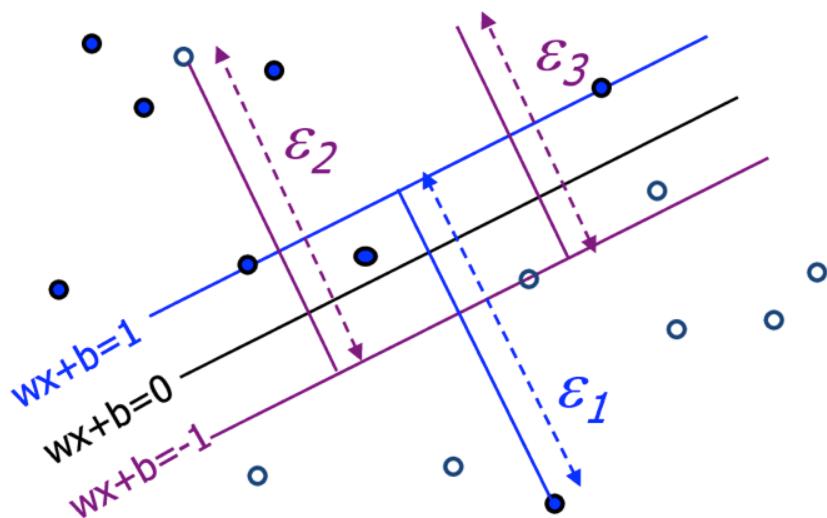
$$f(x) = \sum \alpha_i y_i x_i^T x + b$$

Soft Margin Classification

Slack variables ξ_i can be added to allow misclassification of difficult or noisy examples.

What should our quadratic optimization criterion be?

Minimize



$$\frac{1}{2} \mathbf{w} \cdot \mathbf{w} + C \sum_{k=1}^R \varepsilon_k$$

Soft Margin

$$\min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

The \mathbf{w} that minimizes...

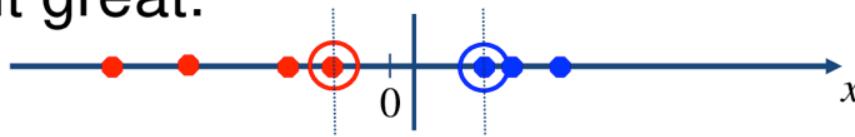
Maximize margin Minimize misclassification

Misclassification cost # data samples Slack variable

subject to $y_i \mathbf{w}^T \mathbf{x}_i \geq 1 - \xi_i,$
 $\xi_i \geq 0, \quad \forall i = 1, \dots, N$

Non-linear SVMs

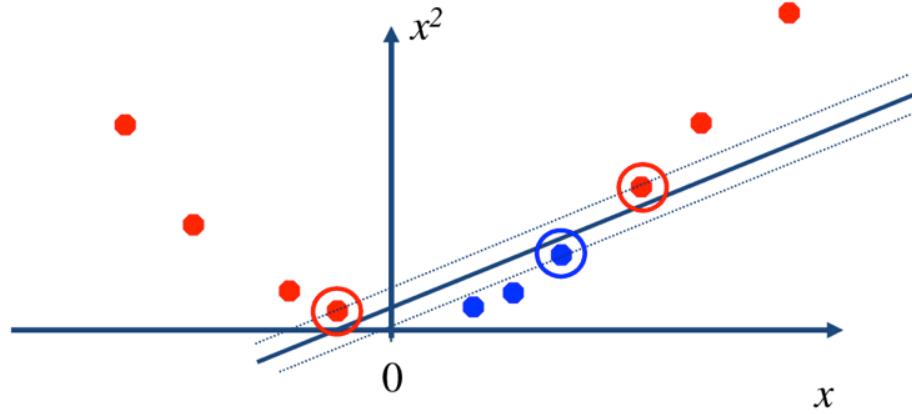
- Datasets that are linearly separable with some noise soft margin work out great:



- But what are we going to do if the dataset is just too hard?



- How about... mapping data to a higher-dimensional space:



The “Kernel Trick”

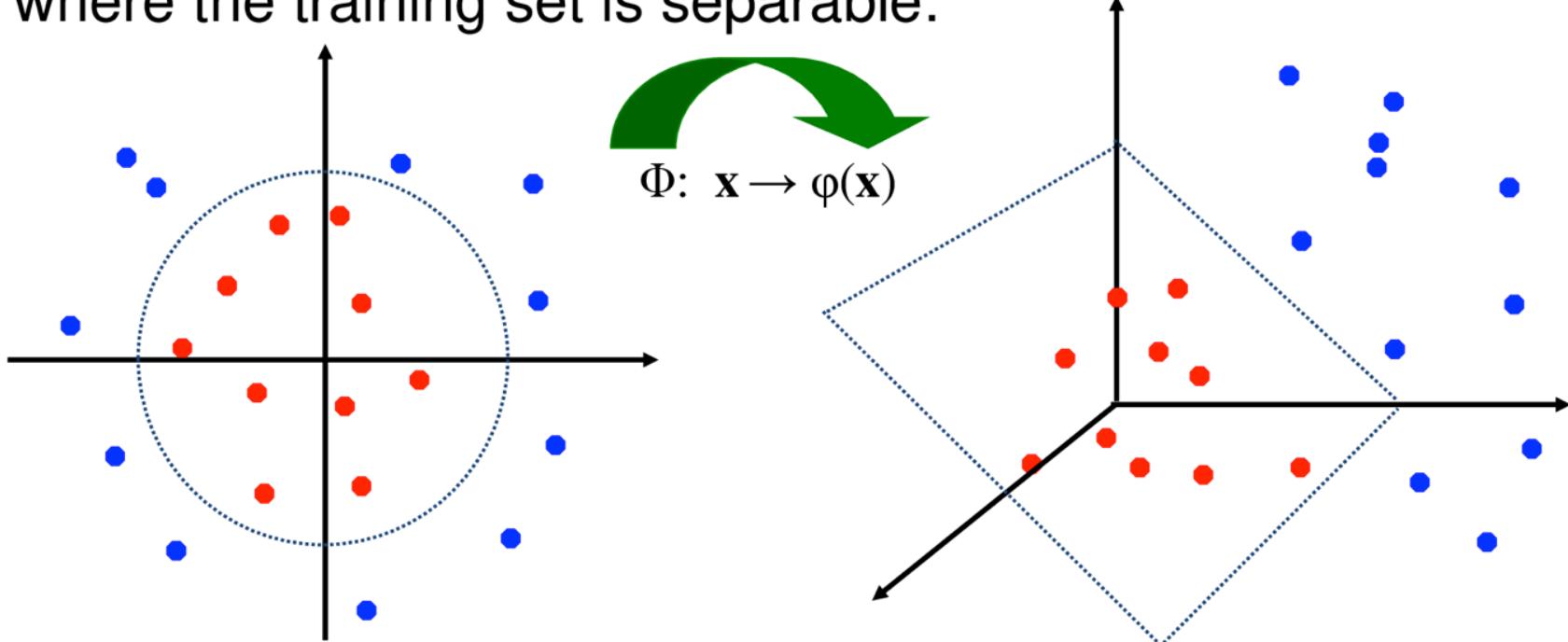
- The linear classifier relies on dot product between vectors
 - $x_i^T \cdot x_j$
- If every data point is mapped into high-dimensional space via some transformation $\Phi: x \rightarrow \phi(x)$, the dot product becomes:
$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$
- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.

SVM Kernels

- SVM algorithms use a set of mathematical functions that are defined as the kernel.
- Function of kernel is to take data as input and transform it into the required form.
- Different SVM algorithms use different types of kernel functions. Example *linear, nonlinear, polynomial, and sigmoid etc.*

Non-linear SVMs: Feature spaces

- General idea: the original input space can always be mapped to some higher-dimensional feature space where the training set is separable:



SVM – Overlapping Class Scenario

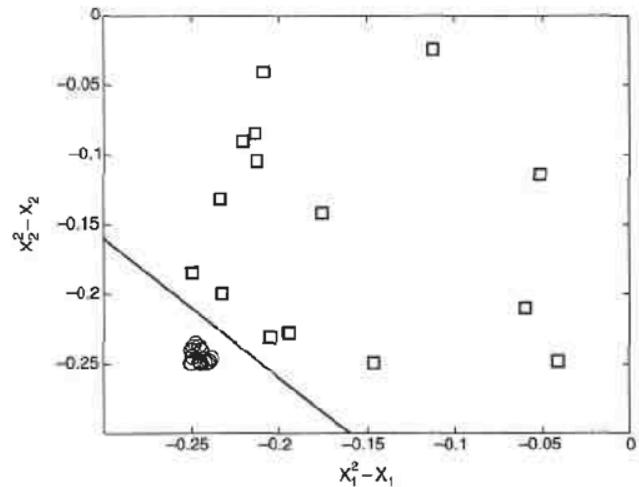
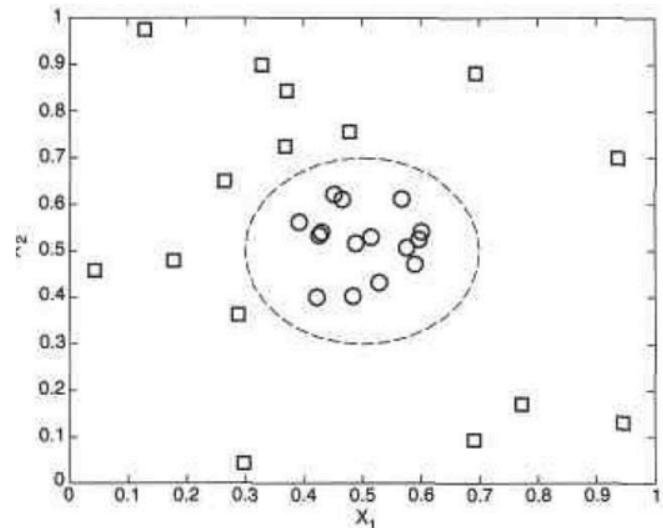
- Data is not separable linearly
- Margin will become inefficient
- Data needs to be transformed from original coordinate space \mathbf{x} to a new space $\Phi(\mathbf{x})$, so that linear decision boundary can be applied
- A non-linear transformation function is needed, like, ex:

$$\Phi : (x_1, x_2) \longrightarrow (x_1^2, x_2^2, \sqrt{2}x_1, \sqrt{2}x_2, 1)$$

- In the transformed space we can choose $\mathbf{w} = (w_0, w_1, \dots, w_4)$ such that

$$w_4x_1^2 + w_3x_2^2 + w_2\sqrt{2}x_1 + w_1\sqrt{2}x_2 + w_0 = 0.$$

- The linear decision boundary in the transformed space has the following form: $\mathbf{w} \cdot \Phi(\mathbf{x}) + b = 0$



Non-linear SVMs Mathematically

- The solution is:

$$f(\mathbf{x}) = \sum \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}_j) + b$$

- Optimization techniques for finding α_i 's remain the same!

Non-linear SVM using kernel

1. Select a kernel function.
 2. Compute pairwise kernel values between labeled examples.
 3. Use this “kernel matrix” to solve for SVM support vectors & alpha weights.
 4. To classify a new example: compute kernel values between new input and support vectors, apply alpha weights, check sign of output.
-

Nonlinear SVM - Overview

- SVM locates a separating hyperplane in the feature space and classify points in that space
 - It does not need to represent the space explicitly, simply by defining a kernel function
 - The kernel function plays the role of the dot product in the feature space.
-

XOR problem: Not linearly separable

For the XOR problem,

\underline{x}_1	$(-1 \quad -1)$	\rightarrow	-1	y_1
\underline{x}_2	$(-1, +1)$	\rightarrow	+1	y_2
\underline{x}_3	$(+1, -1)$	\rightarrow	+1	y_3
\underline{x}_4	$(+1, +1)$	\rightarrow	-1	y_4

XOR problem: Choosing Kernel

Let us consider the XOR problem using SVMs
(Cherkassy, 1998)

Start with a kernel

$$k(\underline{x}, \underline{x}_i) = \left(1 + \underline{x}^T \underline{x}_i \right)^2$$
$$\underline{x} = [x_1 \ x_2]^T \quad \underline{x}_i = [x_{i1} \ x_{i2}]^T$$

$$\begin{aligned}
 k(\underline{x}, \underline{x}_i) &= \left(1 + (\underline{x}_1, \underline{x}_2) \cdot \begin{bmatrix} x_{i1} \\ x_{i2} \end{bmatrix} \right)^2 \\
 &= \left(1 + x_1 x_{i1} + x_2 x_{i2} \right)^2 \\
 &= 1 \cdot 1 + x_1^2 x_{i1}^2 + x_2^2 x_{i2}^2 + 2 x_1 x_2 x_{i1} x_{i2} \\
 &\quad + 2 x_1 x_{i1} + 2 x_2 x_{i2}
 \end{aligned}$$

1 x_{i1}² x_{i2}² 2 x_{i1} x_{i2}
x_{i1}² x_{i2}² 2 x_{i1} 2 x_{i2}

Let us express $k(\cdot, \cdot)$ = $\langle \phi(\underline{x}), \phi(\underline{x}_i) \rangle$

$$\begin{aligned}
 \phi(\underline{x}) &= \begin{bmatrix} 1 & x_1^2 & \sqrt{2} x_1 x_2 & x_2^2 & \sqrt{2} x_1 & \sqrt{2} x_2 \end{bmatrix}^T \\
 \phi(\underline{x}_i) &= \begin{bmatrix} 1 & x_{i1}^2 & \sqrt{2} x_{i1} x_{i2} & x_{i2}^2 & \sqrt{2} x_{i1} & \sqrt{2} x_{i2} \end{bmatrix}^T
 \end{aligned}$$

For the XOR problem,

$$\begin{pmatrix} -1 & -1 \end{pmatrix} \rightarrow -1$$

$$\begin{pmatrix} -1 & +1 \end{pmatrix} \rightarrow +1$$

$$\begin{pmatrix} +1 & -1 \end{pmatrix} \rightarrow +1$$

$$\begin{pmatrix} +1 & +1 \end{pmatrix} \rightarrow -1$$

$$K := \left[K(\underline{x}_i, \underline{x}_j) \right]$$

$$K(\underline{x}_i, \underline{x}_j) = \phi^T(\underline{x}_i) \phi(\underline{x}_j)$$

Each $\underline{x}_i, \underline{x}_j$

$$K = \begin{bmatrix} g & 1 & 1 & 1 \\ 1 & g & 1 & 1 \\ 1 & 1 & g & 1 \\ 1 & 1 & 1 & g \end{bmatrix}$$

Dual of SVM

From the dual problem, the objective $Q(\underline{\alpha})$

$$Q(\underline{\alpha}) = \underline{\alpha}_1 + \underline{\alpha}_2 + \underline{\alpha}_3 + \underline{\alpha}_4 - \frac{1}{2} \left(9\underline{\alpha}_1^2 - 2\underline{\alpha}_1 \underline{\alpha}_2 - 2\underline{\alpha}_1 \underline{\alpha}_3 + 2\underline{\alpha}_1 \underline{\alpha}_4 + 9\underline{\alpha}_2^2 + 2\underline{\alpha}_2 \underline{\alpha}_3 - 2\underline{\alpha}_2 \underline{\alpha}_4 + 9\underline{\alpha}_3^2 - 2\underline{\alpha}_3 \underline{\alpha}_4 + 9\underline{\alpha}_4^2 \right)$$

$\sum_{i=1}^4 \sum_{j=1}^4 \underline{\alpha}_i \underline{\alpha}_j y_i y_j - K(x_i, x_j)$

Solving for α

$$\frac{\partial Q(\underline{\alpha})}{\partial \alpha_i} = 0 \quad i = 1, \dots, 4$$

Doing this yields

$$\left\{ \begin{array}{l} 9\alpha_1 - \alpha_2 - \alpha_3 + \alpha_4 = 1 \\ -\alpha_1 + 9\alpha_2 + \alpha_3 - \alpha_4 = 1 \\ -\alpha_1 + \alpha_2 + 9\alpha_3 - \alpha_4 = 1 \\ \alpha_1 - \alpha_2 - \alpha_3 + 9\alpha_4 = 1 \end{array} \right.$$

Solving this set
of eqns

$$\alpha_{opt,i} = \frac{1}{8} \quad i = 1, \dots, 4$$

Now, all 4 inputs $\{x_i\}_{i=1}^4$ are support vectors

$$Q(\underline{\omega}) = \frac{1}{4}$$

$$\frac{1}{2} \|\underline{\omega}_0\|^2 = \frac{1}{4} \Rightarrow \|\underline{\omega}_0\| = \frac{1}{\sqrt{2}}$$

$$\begin{aligned}\underline{\omega}_0 &= \sum_{i=1}^4 x_{:,i} y_i \phi(x_i) \\ &= \frac{1}{8} \left[-\phi(x_1) + \phi(x_2) + \phi(x_3) - \phi(x_4) \right]\end{aligned}$$

$$\underline{\omega}_0 = [0 \ 0 \ -\frac{1}{\sqrt{2}} \ 0 \ 0 \ 0]^T$$

For +ve support vectors we have

$$\sum_{i=1}^4 \alpha_i y_i K(x_i, x_2) + b = +1 \quad (x_2 \text{ is +ve support vector})$$

Put $x_2 = x_1$ or x_4

$$\alpha_i = \frac{1}{8} \quad \forall i = 1 \text{ to } 4$$

$$K = \begin{bmatrix} 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 1 \end{bmatrix}$$

$$\frac{-K(x_1, x_2)}{8} + \frac{K(x_2, x_2)}{8} + \frac{K(x_3, x_2)}{8} - \frac{K(x_4, x_2)}{8} + b = 1$$

$$-\frac{1}{8} + \frac{1}{8} + \frac{1}{8} - \frac{1}{8} + b = 1$$

$$\boxed{b = 0}$$

Decision boundary of non linear SVM

$$w^T \phi(x) + b = 0$$

The opt. hyperplane is given by

$$\frac{w^T}{w_0} \phi(x) = 0$$

$$\left[0 \quad 0 \quad -\frac{1}{\sqrt{2}} \quad 0 \quad 0 \quad 0 \right] \begin{bmatrix} x_1^2 \\ \sqrt{2}x_1x_2 \\ x_2^2 \\ \sqrt{2}x_1 \\ \sqrt{2}x_2 \end{bmatrix} = 0$$

$$\Rightarrow x_1 x_2 = 0 \text{ is the decision boundary}$$

Non Linear Decision boundary

$$y = -x_1 x_2$$

x_1	x_2	y
-1	-1	-1
-1	+1	+1
+1	-1	+1
+1	+1	-1