

CSE 176 Introduction to Machine Learning Lecture 16: Support Vector Machines

Materials from Olga Vekler and Miguel Carreira-Perpiñán

Support Vector Machine

□Start in 1979 with Vladimir Vapnik's paper

□ Major developments throughout 1990's

SVM v.s. Neural Network
 SVM is better with limited data
 SVM is more interpretable
 SVM is more robust to outlier
 SVM has elegant theory





Topics today

Intuition

Formulation

□ Optimization

□ From linear to non-linear SVM (Kernel SVM)





Intuition of SVM

Linear Discriminant Function

□Which separating hyperplane should we choose?



Margin Intuition

• If sample is close to sample \mathbf{x}_i , it is likely to be on the wrong side



• Poor generalization

Margin Intuition

• Hyperplane as far as possible from any sample



Good generalization

SVM

• Idea: maximize distance to the closest example





SVM: Linearly Separable Case

• SVM: maximize the *margin*



• *margin* is twice the absolute value of distance **b** of the closest example to the separating hyperplane

SVM: Linearly Separable Case



• Support vectors are samples closest to separating hyperplane



Linear SVM

SVM: Formula for the Margin

- $\mathbf{g}(\mathbf{x}) = \mathbf{w}^{\mathsf{t}}\mathbf{x} + \mathbf{w}_{\mathbf{0}}$
- absolute distance between x and the boundary g(x) = 0

$$\frac{\left\|\boldsymbol{W}^{t}\boldsymbol{X}+\boldsymbol{W}_{0}\right\|}{\left\|\boldsymbol{W}\right\|}$$



SVM: Formula for the Margin

- $\mathbf{g}(\mathbf{x}) = \mathbf{w}^{\mathsf{t}}\mathbf{x} + \mathbf{w}_{\mathbf{0}}$
- absolute distance between x and the boundary g(x) = 0

 $\frac{\left\|\boldsymbol{W}^{t}\boldsymbol{X}+\boldsymbol{W}_{0}\right\|}{\left\|\boldsymbol{W}\right\|}$



• distance is unchanged for hyperplane $g_1(x) = \alpha g(x)$

$$\frac{\left|\alpha \mathbf{w}^{\mathsf{t}} \mathbf{x} + \alpha \mathbf{w}_{0}\right|}{\left\|\alpha \mathbf{w}\right\|} = \frac{\left|\mathbf{w}^{\mathsf{t}} \mathbf{x} + \mathbf{w}_{0}\right|}{\left\|\mathbf{w}\right\|}$$

Let x_i be an example closest to the boundary. Set

$$|\mathbf{w}^{\mathsf{t}}\mathbf{X}_{\mathsf{i}} + \mathbf{w}_{\mathsf{0}}| = 1$$

• Now the largest margin hyperplane is unique

SVM: Formula for the Margin

- For uniqueness, set $|\mathbf{w}^{t}\mathbf{x}_{i} + \mathbf{w}_{0}| = 1$ for any example \mathbf{x}_{i} closest to the boundary
- now distance from closest sample \mathbf{x}_i to $\mathbf{g}(\mathbf{x}) = 0$ is

$$\frac{\mathbf{w}^{\mathsf{t}}\mathbf{x}_{i}+\mathbf{w}_{0}}{\|\mathbf{w}\|}=\frac{1}{\|\mathbf{w}\|}$$

• Thus the margin is

$$\mathbf{m} = \frac{2}{\|\mathbf{w}\|}$$



SVM: Optimal Hyperplane

- Maximize margin
 - subject to constraints
 - $\begin{cases} \mathbf{w}^{\mathsf{t}} \mathbf{x}_{i} + \mathbf{w}_{0} \ge 1 & \text{if } \mathbf{x}_{i} \text{ is positive example} \\ \mathbf{w}^{\mathsf{t}} \mathbf{x}_{i} + \mathbf{w}_{0} \le -1 & \text{if } \mathbf{x}_{i} \text{ is negative example} \end{cases}$
- Let $\begin{cases} z_i = 1 & \text{if } x_i \text{ is positive example} \\ z_i = -1 & \text{if } x_i \text{ is negative example} \end{cases}$
- Convert our problem to

minimize
$$J(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2$$

constrained to $\mathbf{z}^i (\mathbf{w}^t \mathbf{x}_i + \mathbf{w}_0) \ge 1 \quad \forall \mathbf{i}$

• J(w) is a convex function, thus it has a single global minimum

$$\mathbf{n} = \frac{2}{\|\mathbf{w}\|}$$

SVM: Optimal Hyperplane

• Use Kuhn-Tucker theorem to convert our problem to:

maximize
$$\mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{j} \mathbf{z}_{i} \mathbf{z}_{j} \mathbf{x}_{i}^{\mathsf{t}} \mathbf{x}_{j}$$
constrained to
$$\alpha_{i} \ge 0 \quad \forall \mathbf{i} \quad \text{and} \quad \sum_{i=1}^{n} \alpha_{i} \mathbf{z}_{i} = 0$$

- $\alpha = {\alpha_1, ..., \alpha_n}$ are new variables, one for each sample
- $L_D(\alpha)$ can be optimized by quadratic programming
- $L_{D}(\alpha)$ formulated in terms of α
 - depends on w and w₀

SVM: Optimal Hyperplane

- After finding the optimal $\alpha = \{\alpha_1, ..., \alpha_n\}$
 - for every sample **i**, one of the following must hold
 - $\alpha_i = 0$ (sample *i* is not a support vector)
 - $\alpha_i \neq 0$ and $\mathbf{z}_i(\mathbf{w}^t \mathbf{x}_i + \mathbf{w}_0 1) = 0$ (sample **i** is support vector)
 - compute $\mathbf{w} = \sum_{i=1}^{n} \alpha_i \mathbf{z}_i \mathbf{x}_i$ • solve for \mathbf{w}_0 using any $\alpha_i > 0$ and $\alpha_i [\mathbf{z}_i (\mathbf{w}^t \mathbf{x}_i + \mathbf{w}_0) - 1] = 0$ $\mathbf{w}_0 = \frac{1}{\mathbf{z}_i} - \mathbf{w}^t \mathbf{x}_i$
- Final discriminant function:

$$\mathbf{g}(\mathbf{x}) = \left(\sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \mathbf{x}_i\right)^{\mathbf{t}} \mathbf{x} + \mathbf{w}_0$$

• where **S** is the set of support vectors

$$\mathbf{S} = \left\{ \mathbf{x}_{i} \mid \boldsymbol{\alpha}_{i} \neq \mathbf{0} \right\}$$

SVM: Non Separable Case

• Linear classifier still be appropriate when data is not linearly separable, but almost linearly separable



• Can adapt SVM to almost linearly separable case

SVM as Unconstrained Minimization

• SVM objective can be rewritten as unconstrained optimization



- z_i f(x_i) > 1 : x_i is on the right side of the hyperplane and outside margin, no loss
- z_i f(x_i) = 1 : x_i on the margin, no loss
- z_i f(x_i) < 1 : x_i is inside margin, or on the wrong side of the hyperplane, contributes to loss

SVM: Hinge Loss

• SVM uses Hinge loss per sample **x**_i

$$\mathbf{L}_{i}(\mathbf{x}_{i}) = \max(0, 1 - \mathbf{z}_{i}\mathbf{f}(\mathbf{x}_{i}))$$



Hinge loss encourages classification with a margin of 1

SVM: Hinge Loss

• Can optimize with gradient descent, convex function

$$J(\mathbf{w}) = \frac{1}{2} \|\mathbf{w}\|^2 + \beta \sum_{i=1}^{n} \max(0, 1 - \mathbf{z}_i \mathbf{f}(\mathbf{x}_i))$$
$$\mathbf{f}(\mathbf{x}_i) = \mathbf{w}^{\mathsf{t}} \mathbf{x}_i + \mathbf{w}_0$$

• Gradient $\mathbf{L}(\mathbf{x}_i)$ $\frac{\partial \mathbf{J}}{\partial \mathbf{w}} = \mathbf{w} - \mathbf{z}_i \mathbf{x}_i$ $\frac{\partial \mathbf{J}}{\partial \mathbf{w}} = \mathbf{w}$ $\mathbf{z}_i \mathbf{f}(\mathbf{x}_i)$

 Gradient descent, single sample

$$\mathbf{w} = \begin{cases} \mathbf{w} - \alpha (\mathbf{w} - \beta \mathbf{z}_i \mathbf{x}_i) & \text{if } \mathbf{z}_i \mathbf{f}(\mathbf{x}_i) < 1 \\ \mathbf{w} - \alpha \mathbf{w} & \text{otherwise} \end{cases}$$



Non-linear SVM (Kernel SVM)

Non Linear Mapping

- Cover's theorem:
 - *"pattern-classification problem cast in a high dimensional space non-linearly is more likely to be linearly separable than in a low-dimensional space"*
- Not linearly separable in 1D
 Lift to 2D space with h(x) = (x,x²)





Non Linear Mapping

- To solve a non linear problem with a linear classifier
 - 1. Project data \mathbf{x} to high dimension using function $\boldsymbol{\varphi}(\mathbf{x})$
 - 2. Find a linear discriminant function for transformed data $\varphi(\mathbf{x})$
 - 3. Final nonlinear discriminant function is $\mathbf{g}(\mathbf{x}) = \mathbf{w}^{t} \boldsymbol{\varphi}(\mathbf{x}) + \mathbf{w}_{0}$



• In 2D, discriminant function is linear

$$\mathbf{g}\!\left(\!\begin{bmatrix}\mathbf{x}^{(1)}\\\mathbf{x}^{(2)}\end{bmatrix}\!\right)\!=\!\begin{bmatrix}\mathbf{w}_1 & \mathbf{w}_2\end{bmatrix}\!\begin{bmatrix}\mathbf{x}^{(1)}\\\mathbf{x}^{(2)}\end{bmatrix}\!+\mathbf{w}_0$$

• In 1D, discriminant function is not linear $\mathbf{g}(\mathbf{x}) = \mathbf{w}_1 \mathbf{x} + \mathbf{w}_2 \mathbf{x}^2 + \mathbf{w}_0$

Non Linear Mapping: Another Example



Non Linear SVM

- Can use any linear classifier after lifting data into a higher dimensional space
- However we will have to deal with the "curse of dimensionality"
 - 1. poor generalization to test data
 - 2. computationally expensive
- SVM avoids the "curse of dimensionality" by
 - enforcing largest margin permits good generalization
 - computation in the higher dimensional case is performed only implicitly through the use of *kernel* functions

Recall SVM optimization

maximize
$$\mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{i} \mathbf{z}_{i} \mathbf{z}_{j} \mathbf{x}_{i}^{\mathsf{t}} \mathbf{x}_{j}$$

- Optimization depends on samples x_i only through the dot product x_i^tx_j
- If we lift x_i to high dimension using $\varphi(\mathbf{x})$, need to compute high dimensional product $\varphi(x_i)^t \varphi(x_j)$

maximize
$$\mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{i} \mathbf{z}_{i} \mathbf{z}_{j} \frac{\varphi(\mathbf{x}_{i})^{t} \varphi(\mathbf{x}_{j})}{\mathbf{K}(\mathbf{x}_{i}, \mathbf{x}_{j})}$$

• Idea: find *kernel* function $K(\mathbf{x}_i, \mathbf{x}_j)$ s.t. $K(\mathbf{x}_i, \mathbf{x}_j) = \boldsymbol{\varphi}(\mathbf{x}_i)^t \boldsymbol{\varphi}(\mathbf{x}_j)$

maximize
$$\mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{i} \mathbf{z}_{i} \mathbf{z}_{j} \boldsymbol{\varphi}(\mathbf{x}_{i})^{t} \boldsymbol{\varphi}(\mathbf{x}_{j})$$
$$\mathbf{K}(\mathbf{x}_{i}, \mathbf{x}_{j})$$

- Kernel trick
 - only need to compute $K(\mathbf{x}_i, \mathbf{x}_i)$ instead of $\boldsymbol{\varphi}(\mathbf{x}_i)^t \boldsymbol{\varphi}(\mathbf{x}_i)$
 - no need to lift data in high dimension explicitly, computation is performed in the original dimension

- Suppose we have 2 features and $K(x,y) = (x^ty)^2$
- Which mapping $\boldsymbol{\varphi}(\mathbf{x})$ does it correspond to?

$$\begin{split} \mathbf{K}(\mathbf{x},\mathbf{y}) &= \left(\mathbf{x}^{t}\mathbf{y}\right)^{2} = \left(\begin{bmatrix} \mathbf{x}^{(1)} & \mathbf{x}^{(2)} \end{bmatrix} \begin{bmatrix} \mathbf{y}^{(1)} \\ \mathbf{y}^{(2)} \end{bmatrix} \right)^{2} = \left(\mathbf{x}^{(1)}\mathbf{y}^{(1)} + \mathbf{x}^{(2)}\mathbf{y}^{(2)}\right)^{2} \\ &= \left(\mathbf{x}^{(1)}\mathbf{y}^{(1)}\right)^{2} + 2\left(\mathbf{x}^{(1)}\mathbf{y}^{(1)}\right)\left(\mathbf{x}^{(2)}\mathbf{y}^{(2)}\right) + \left(\mathbf{x}^{(2)}\mathbf{y}^{(2)}\right)^{2} \\ &= \begin{bmatrix} \left(\mathbf{x}^{(1)}\right)^{2} & \sqrt{2}\mathbf{x}^{(1)}\mathbf{x}^{(1)}\mathbf{x}^{(2)} & \left(\mathbf{x}^{(2)}\right)^{2} \end{bmatrix} \begin{bmatrix} \left(\mathbf{y}^{(1)}\right)^{2} & \sqrt{2}\mathbf{y}^{(1)}\mathbf{y}^{(1)}\mathbf{y}^{(2)} & \left(\mathbf{y}^{(2)}\right)^{2} \end{bmatrix}^{t} \end{split}$$

Thus

$$\boldsymbol{\phi}(\mathbf{x}) = \begin{bmatrix} \mathbf{x}^{(1)} \\ \mathbf{x}^{(2)} \end{bmatrix} \begin{bmatrix} \mathbf{x}^{(2)} \\ \mathbf{x}^{(2)} \end{bmatrix}$$



- How to choose kernel K(x_i,x_i)?
 - $K(\mathbf{x}_i, \mathbf{x}_j)$ should correspond to product $\boldsymbol{\varphi}(\mathbf{x}_i)^t \boldsymbol{\varphi}(\mathbf{x}_j)$ in a higher dimensional space
 - Mercer's condition states which kernel function can be expressed as dot product of two vectors
 - Kernel's not satisfying Mercer's condition can be sometimes used, but no geometrical interpretation
- Common choices satisfying Mercer's condition
 - Polynomial kernel $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^{t} \mathbf{x}_j + 1)^{p}$
 - Gaussian radial Basis kernel (data is lifted in infinite dimensions)

$$\mathbf{K}(\mathbf{x}_{i},\mathbf{x}_{j}) = \exp\left(-\frac{1}{2\sigma^{2}} \|\mathbf{x}_{i}-\mathbf{x}_{j}\|^{2}\right)$$

Non Linear SVM

search for separating hyperplane in high dimension

$$\mathbf{w}\boldsymbol{\phi}(\mathbf{x}) + \mathbf{w}_0 = \mathbf{0}$$

Choose φ(x) so that the first ("0"th) dimension is the augmented dimension with feature value fixed to 1

$$\boldsymbol{\varphi}(\mathbf{x}) = \begin{bmatrix} 1 & \mathbf{x}^{(1)} & \mathbf{x}^{(2)} & \mathbf{x}^{(1)}\mathbf{x}^{(2)} \end{bmatrix}^{\mathbf{t}}$$

• Threshold \mathbf{w}_0 gets folded into vector \mathbf{w}

$$\begin{bmatrix} \mathbf{w}_0 & \mathbf{w} \end{bmatrix} \begin{bmatrix} \mathbf{1} \\ \mathbf{*} \end{bmatrix} = \mathbf{0}$$
$$\mathbf{\phi}(\mathbf{x})$$

Non Linear SVM

• Thus seeking hyperplane

$$\mathbf{w}\phi(\mathbf{x}) = \mathbf{0}$$

- Or, equivalently, a hyperplane that goes through the origin in high dimensions
 - removes only one degree of freedom
 - but we introduced many new degrees when lifted the data in high dimension

Non Linear SVM Recepie

- Start with $\mathbf{x}_1, \dots, \mathbf{x}_n$ in original feature space of dimension **d**
- Choose kernel **K**(**x**_i,**x**_j)
 - implicitly chooses function $\boldsymbol{\varphi}(\mathbf{x}_i)$ that takes \mathbf{x}_i to a higher dimensional space
 - gives dot product in the high dimensional space
- Find largest margin linear classifier in the higher dimensional space by using quadratic programming package to solve

$$\begin{array}{ll} \text{maximize} & \mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{n} \alpha_{i} - \frac{1}{2} \sum_{i=1}^{n} \sum_{j=1}^{n} \alpha_{i} \alpha_{i} \mathbf{z}_{i} \mathbf{z}_{j} \mathbf{K}(\mathbf{x}_{i}, \mathbf{x}_{j}) \\ \\ \text{constrained to} & 0 \leq \alpha_{i} \leq \beta \quad \forall i \quad \text{and} \ \sum_{i=1}^{n} \alpha_{i} \mathbf{z}_{i} = 0 \end{array}$$

Non Linear SVM Recipe

• Weight vector **w** in the high dimensional space

$$\mathbf{w} = \sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \boldsymbol{\phi}(\mathbf{x}_i)$$

• where **S** is the set of support vectors

$$\mathbf{S} = \left\{ \mathbf{x}_{i} \mid \boldsymbol{\alpha}_{i} \neq \mathbf{0} \right\}$$

- Linear discriminant function in the high dimensional space $\mathbf{g}(\boldsymbol{\varphi}(\mathbf{x})) = \mathbf{w}^{\mathsf{t}} \boldsymbol{\varphi}(\mathbf{x}) = \left(\sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \boldsymbol{\varphi}(\mathbf{x}_i)\right)^{\mathsf{t}} \boldsymbol{\varphi}(\mathbf{x})$
- Non linear discriminant function in the original space:

$$\mathbf{g}(\mathbf{x}) = \left(\sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \boldsymbol{\phi}(\mathbf{x}_i)\right)^{\mathsf{t}} \boldsymbol{\phi}(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \boldsymbol{\phi}^{\mathsf{t}}(\mathbf{x}_i) \boldsymbol{\phi}(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \mathsf{K}(\mathbf{x}_i, \mathbf{x})$$

• Decide class 1 if g(x) > 0, otherwise decide class 2

Non Linear SVM

• Nonlinear discriminant function

$$g(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathbf{S}} \alpha_i \mathbf{z}_i \mathbf{K}(\mathbf{x}_i, \mathbf{x})$$

$$g(\mathbf{x}) = \sum_{\mathbf{x}_i \in \mathbf{S}} \text{weight of support}_{vector \ \mathbf{x}_i} \mathbf{F1} \text{similarity}_{between \ \mathbf{x} \text{ and}_{support vector \ \mathbf{x}_i}}$$

$$f(\mathbf{x}_i, \mathbf{x}) = \exp\left(-\frac{1}{2\sigma^2} \|\mathbf{x}_i - \mathbf{x}\|^2\right)$$

- Class 1: **x**₁ = [1,-1], **x**₂ = [-1,1]
- Class 2: **x**₃ = [1,1], **x**₄ = [-1,-1]
- Use polynomial kernel of degree 2

•
$$K(x_i, x_j) = (x_i^{t} x_j + 1)^2$$

constrained to

• Kernel corresponds to mapping $\left(\mathbf{x}\right) = \begin{bmatrix} 1 & \sqrt{2}\mathbf{x}^{(1)} & \sqrt{2}\mathbf{x}^{(2)} & \sqrt{2}\mathbf{x}^{(1)}\mathbf{x}^{(2)} & \left(\mathbf{x}^{(1)}\right)^2 & \left(\mathbf{x}^{(2)}\right)^2 \end{bmatrix}^t$

• Need to maximize
$$\mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{4} \alpha_{i} \quad \frac{1}{2} \sum_{i=1}^{4} \sum_{j=1}^{4} \alpha_{j} \alpha_{i} \alpha_{j} \mathbf{z}_{i} \mathbf{z}_{j} (\mathbf{x}_{i}^{\mathsf{t}} \mathbf{x}_{j} + 1)^{2}$$

 $0 \leq \alpha_{i} \hspace{0.1in} \forall i \hspace{0.1in} \text{and} \hspace{0.1in} \alpha_{_{1}} + \alpha_{_{2}} - \alpha_{_{3}} - \alpha_{_{4}} = 0$

OO

- Rewrite $\mathbf{L}_{\mathbf{D}}(\alpha) = \sum_{i=1}^{4} \alpha_i \frac{1}{2} \alpha^t \mathbf{H} \alpha$ • where $\alpha = [\alpha_1 \quad \alpha_2 \quad \alpha_3 \quad \alpha_4]^t$ and $\mathbf{H} = \begin{bmatrix} 9 & 1 & -1 & -1 \\ 1 & 9 & -1 & -1 \\ -1 & -1 & 9 & 1 \\ -1 & -1 & 1 & 9 \end{bmatrix}$
- Take derivative with respect to α and set it to $\boldsymbol{0}$

$$\frac{\mathbf{d}}{\mathbf{da}}\mathbf{L}_{\mathbf{D}}(\alpha) = \begin{bmatrix} 1\\1\\1\\1\\1 \end{bmatrix} - \begin{bmatrix} 9 & 1 & -1 & -1\\1 & 9 & -1 & -1\\-1 & -1 & 9 & 1\\-1 & -1 & 1 & 9 \end{bmatrix} \alpha = 0$$

- Solution to the above is $\alpha_1 = \alpha_2 = \alpha_3 = \alpha_4 = 0.25$
 - satisfies the constraints $\forall i$, $0 \le \alpha_i$ and $\alpha_1 + \alpha_2 \alpha_3 \alpha_4 = 0$
 - all samples are support vectors

$$(\mathbf{x}) = \begin{bmatrix} 1 & \sqrt{2}\mathbf{x}^{(1)} & \sqrt{2}\mathbf{x}^{(2)} & \sqrt{2}\mathbf{x}^{(1)}\mathbf{x}^{(2)} & (\mathbf{x}^{(1)})^2 & (\mathbf{x}^{(2)})^2 \end{bmatrix}^{t}$$

• Weight vector **w** is:

$$\mathbf{w} = \sum_{i=1}^{4} \alpha_{i} \mathbf{z}_{i} \phi(\mathbf{x}_{i}) = 0.25(\phi(\mathbf{x}_{1}) + \phi(\mathbf{x}_{2}) - \phi(\mathbf{x}_{3}) - \phi(\mathbf{x}_{4}))$$
$$= \begin{bmatrix} 0 & 0 & 0 & \sqrt{2} & 0 & 0 \end{bmatrix}$$

• by plugging in $\mathbf{x_1} = [1,-1], \mathbf{x_2} = [-1,1], \mathbf{x_3} = [1,1], \mathbf{x_4} = [-1,-1]$

• Nonlinear discriminant function is

$$g(\mathbf{x}) = \mathbf{w} \phi(\mathbf{x}) = \sum_{i=1}^{6} \mathbf{w}_{i} \phi_{i}(\mathbf{x}) = \sqrt{2} \left(\sqrt{2} \mathbf{x}^{(1)} \mathbf{x}^{(2)} \right) = 2 \mathbf{x}^{(1)} \mathbf{x}^{(2)}$$



decision boundaries nonlinear

decision boundary is linear

Degree 3 Polynomial Kernel



- Left: In linearly separable case, decision boundary is roughly linear, indicating that dimensionality is controlled
- Right: nonseparable case is handled by a polynomial of degree 3

SVM Summary

- Advantages:
 - nice theory
 - good generalization properties
 - objective function has no local minima
 - can be used to find non linear discriminant functions
 - often works well in practice, even if not a lot of training data
- Disadvantages:
 - tends to be slower than other methods
 - quadratic programming is computationally expensive