

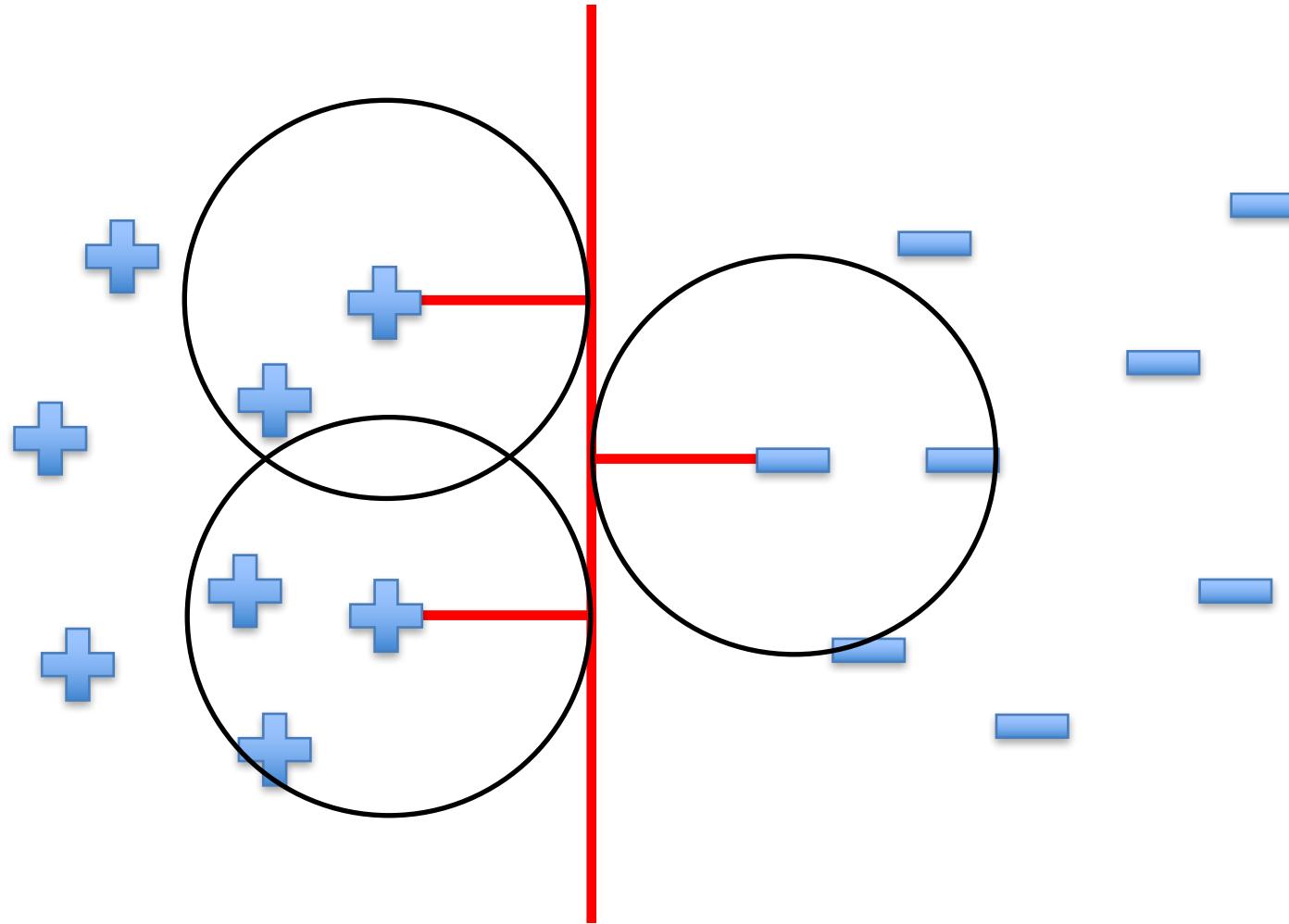


# Support Vector Machines & Kernels

Doing *really* well with linear decision surfaces

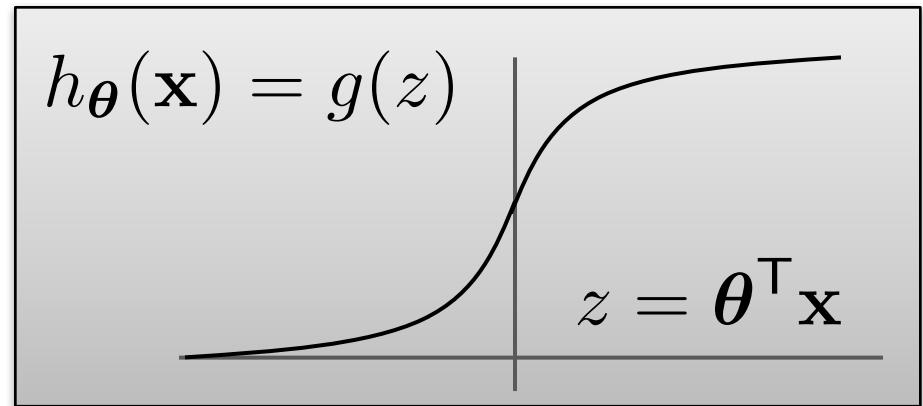
These slides were assembled by Eric Eaton, with grateful acknowledgement of the many others who made their course materials freely available online. Feel free to reuse or adapt these slides for your own academic purposes, provided that you include proper attribution. Please send comments and corrections to Eric.

# Maximizing the Margin



# Alternative View of Logistic Regression

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$



If  $y = 1$ , we want  $h_{\theta}(\mathbf{x}) \approx 1$ ,  $\theta^T \mathbf{x} \gg 0$

If  $y = 0$ , we want  $h_{\theta}(\mathbf{x}) \approx 0$ ,  $\theta^T \mathbf{x} \ll 0$

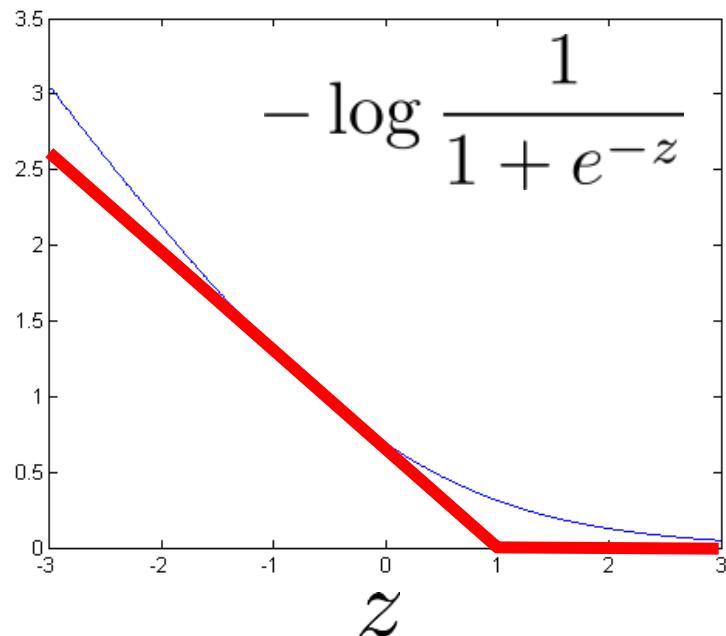
$$J(\theta) = - \sum_{i=1}^n [y_i \log h_{\theta}(\mathbf{x}_i) + (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))]$$
$$\min_{\theta} J(\theta) \quad \underbrace{\text{cost}_1(\theta^T \mathbf{x}_i)} \quad \underbrace{\text{cost}_0(\theta^T \mathbf{x}_i)}$$

# Alternate View of Logistic Regression

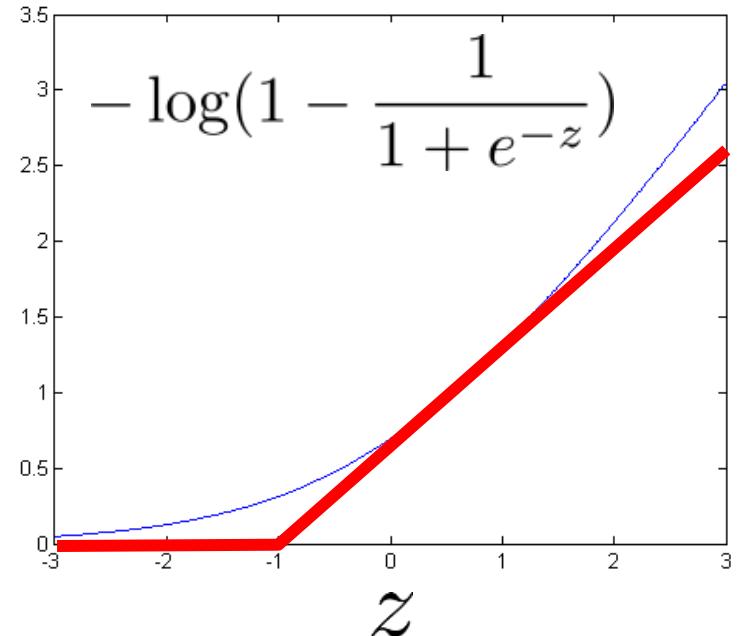
Cost of example:  $-y_i \log h_{\theta}(\mathbf{x}_i) - (1 - y_i) \log (1 - h_{\theta}(\mathbf{x}_i))$

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}} \quad z = \theta^T \mathbf{x}$$

If  $y = 1$  (want  $\theta^T \mathbf{x} \gg 0$ ):

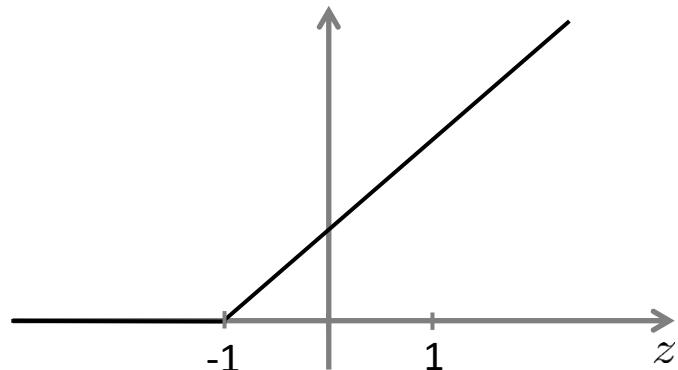
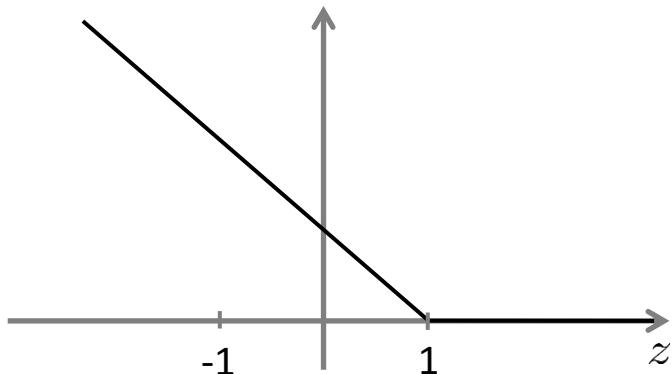


If  $y = 0$  (want  $\theta^T \mathbf{x} \ll 0$ ):



# Support Vector Machine

If  $y = 1$  (want  $\theta^\top \mathbf{x} \geq 1$ ):      If  $y_i = -1$  (want  $\theta^\top \mathbf{x} \leq -1$ ):



$$\ell_{\text{hinge}}(h(\mathbf{x})) = \max(0, 1 - y \cdot h(\mathbf{x}))$$

# Support Vector Machine

$$\min_{\theta} C \sum_{i=1}^n [y_i \text{cost}_1(\theta^\top \mathbf{x}_i) + (1 - y_i) \text{cost}_0(\theta^\top \mathbf{x}_i)] + \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

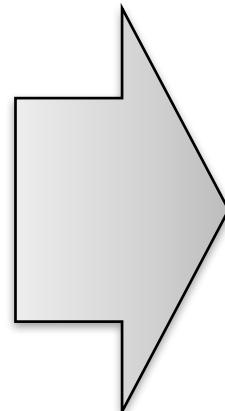
$y = 1 / 0$

with  $C = 1$

$y = +1 / -1$

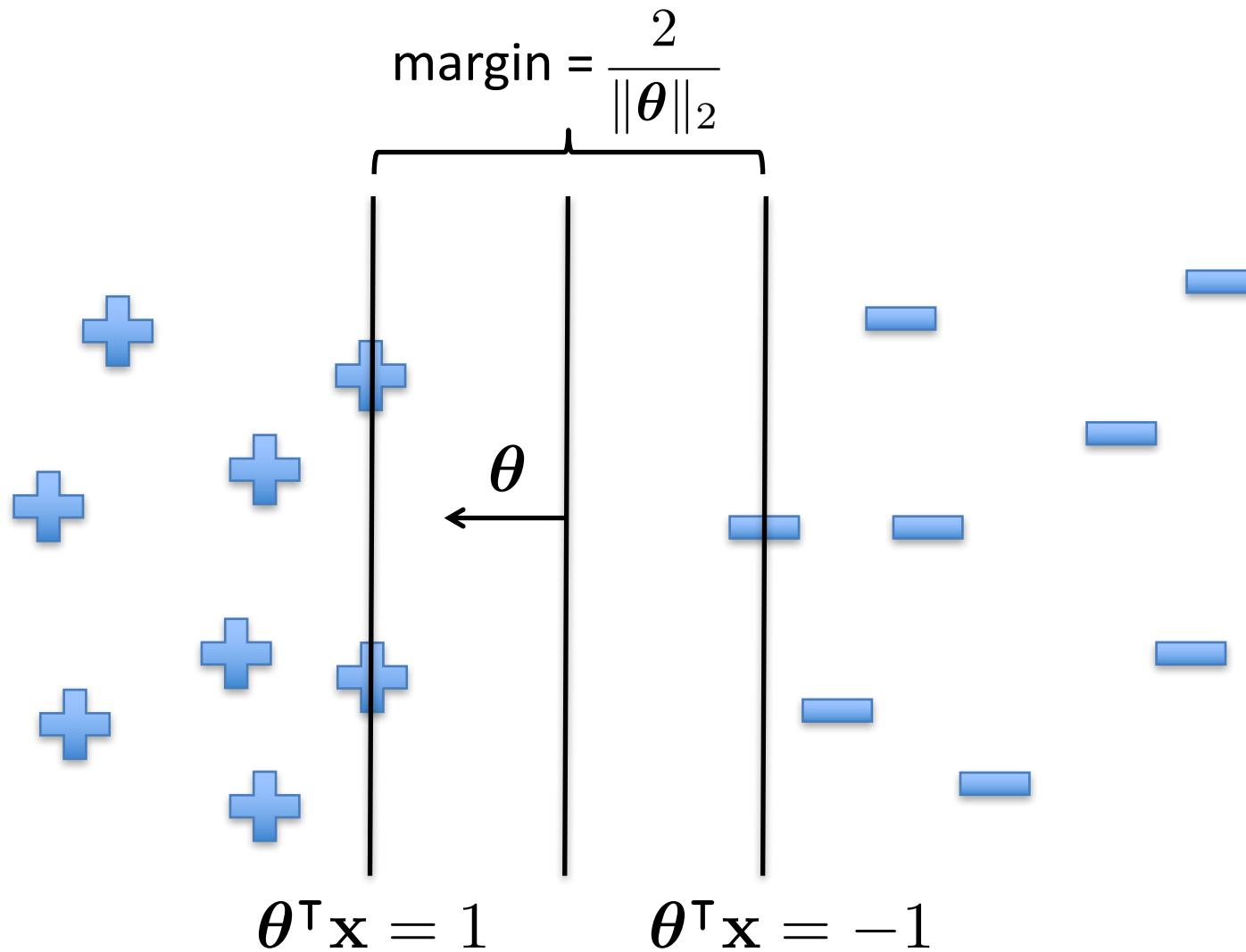
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

$$\begin{aligned} \text{s.t. } \theta^\top \mathbf{x}_i &\geq 1 & \text{if } y_i = 1 \\ \theta^\top \mathbf{x}_i &\leq -1 & \text{if } y_i = -1 \end{aligned}$$

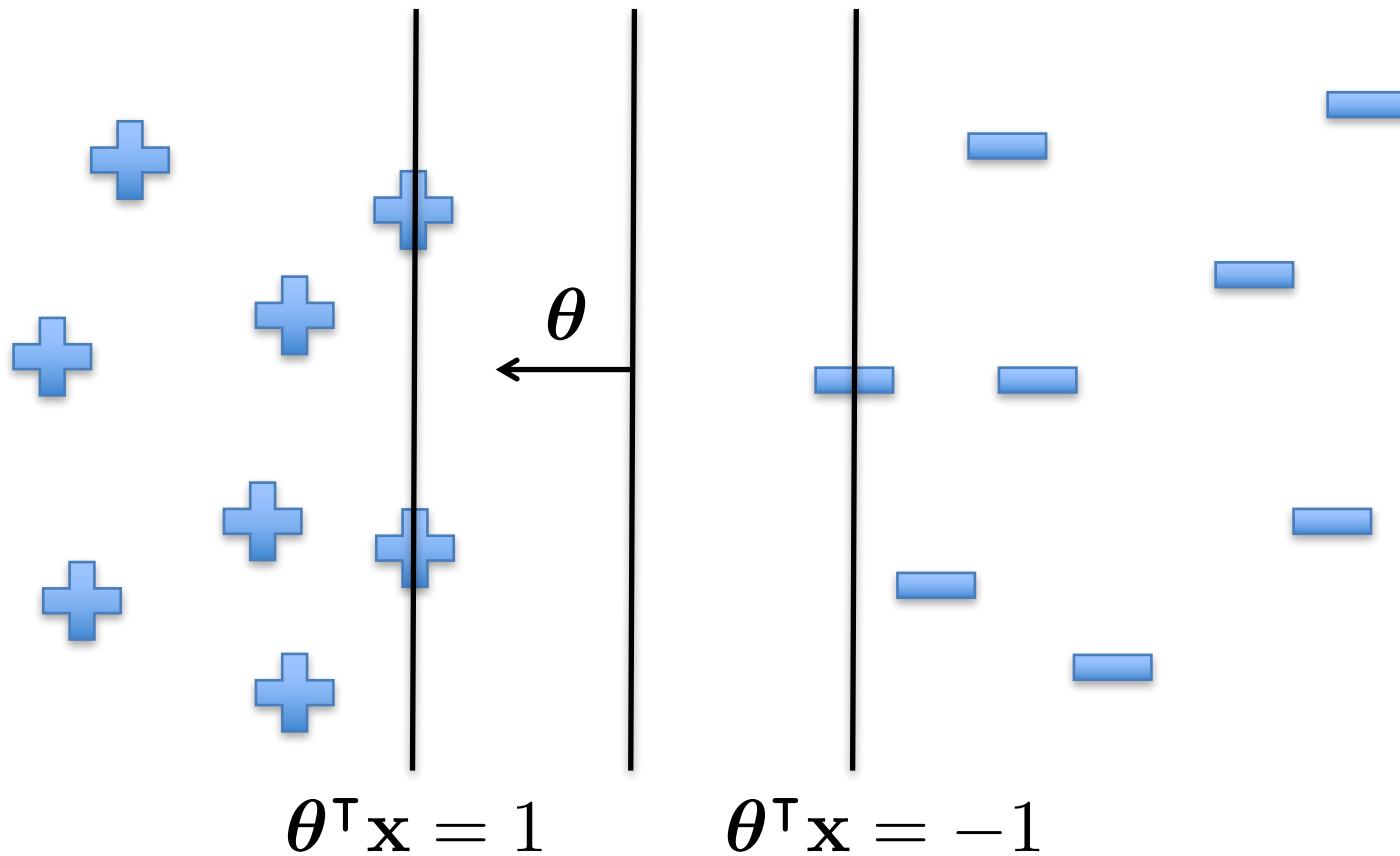


$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$
$$\text{s.t. } y_i(\theta^\top \mathbf{x}_i) \geq 1$$

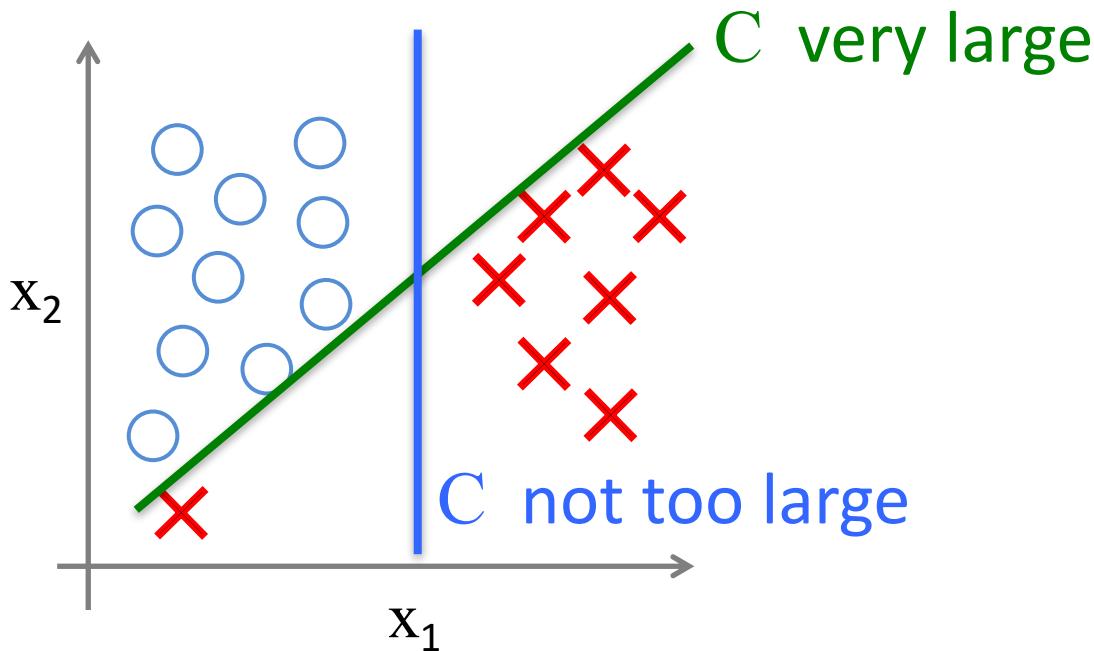
# Maximum Margin Hyperplane



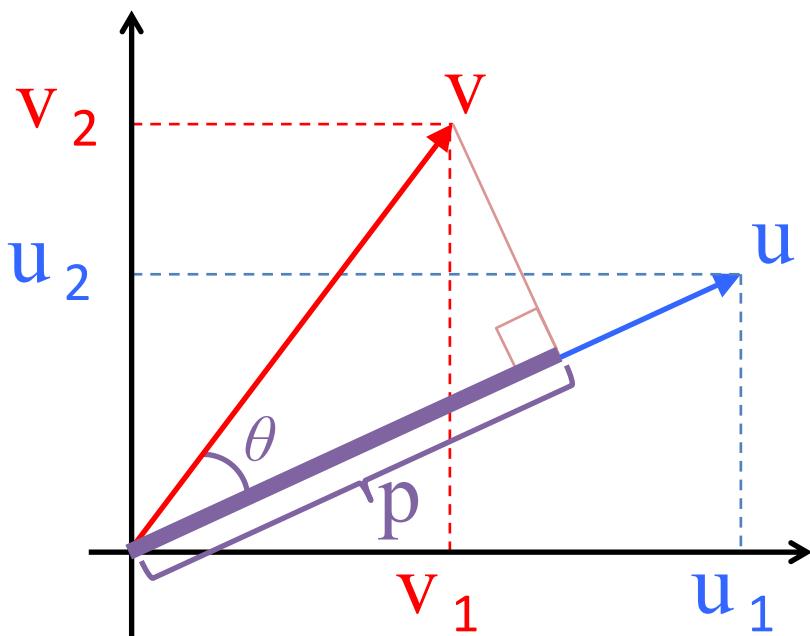
# Support Vectors



# Large Margin Classifier in Presence of Outliers



# Vector Inner Product



$$u = \begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \quad v = \begin{bmatrix} v_1 \\ v_2 \end{bmatrix}$$

$$\begin{aligned}\|\mathbf{u}\|_2 &= \text{length}(\mathbf{u}) \in \mathbb{R} \\ &= \sqrt{u_1^2 + u_2^2}\end{aligned}$$

$$\begin{aligned}\mathbf{u}^\top \mathbf{v} &= \mathbf{v}^\top \mathbf{u} \\ &= u_1 v_1 + u_2 v_2 \\ &= \|\mathbf{u}\|_2 \|\mathbf{v}\|_2 \cos \theta \\ &= p \|\mathbf{u}\|_2 \quad \text{where } p = \|\mathbf{v}\|_2 \cos \theta\end{aligned}$$

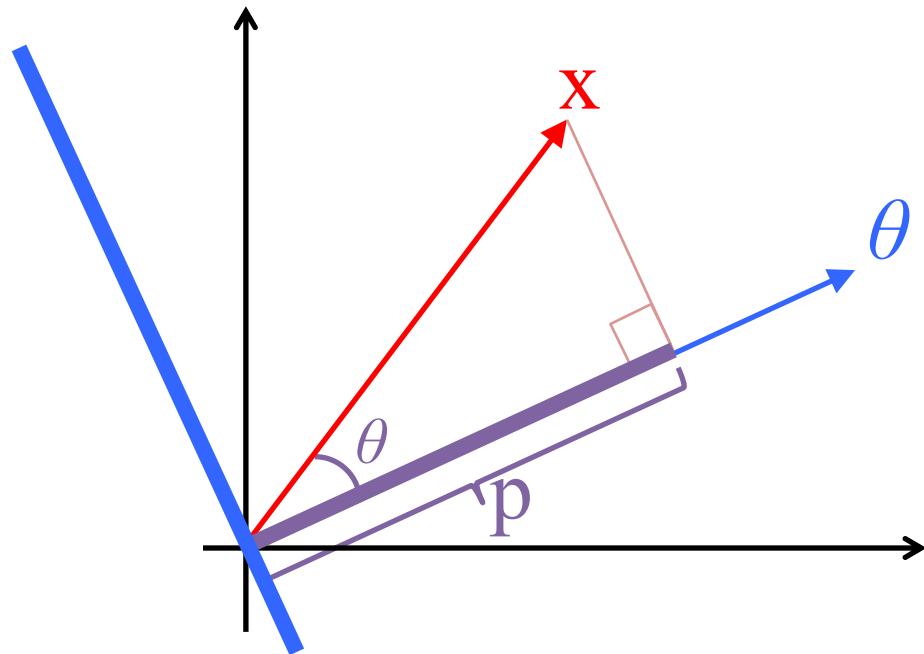
# Understanding the Hyperplane

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

s.t.  $\theta^\top \mathbf{x}_i \geq 1 \quad \text{if } y_i = 1$

$\theta^\top \mathbf{x}_i \leq -1 \quad \text{if } y_i = -1$

Assume  $\theta_0 = 0$  so that the hyperplane is centered at the origin, and that  $d = 2$



$$\begin{aligned}\theta^\top \mathbf{x} &= \|\theta\|_2 \underbrace{\|\mathbf{x}\|_2 \cos \theta}_p \\ &= p \|\theta\|_2\end{aligned}$$

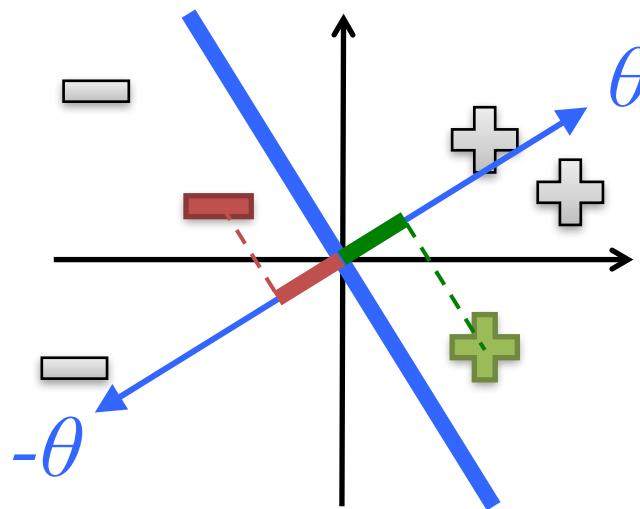
# Maximizing the Margin

$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2$$

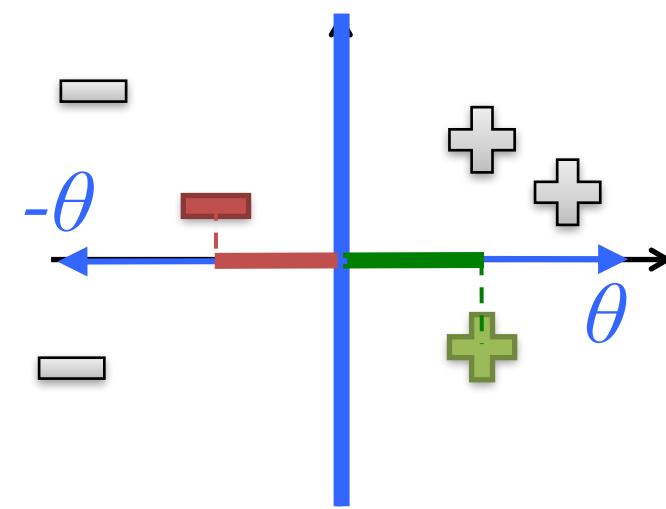
s.t.  $\theta^\top \mathbf{x}_i \geq 1 \quad \text{if } y_i = 1$   
 $\theta^\top \mathbf{x}_i \leq -1 \quad \text{if } y_i = -1$

Assume  $\theta_0 = 0$  so that the hyperplane is centered at the origin, and that  $d = 2$

Let  $p_i$  be the projection of  $\mathbf{x}_i$  onto the vector  $\theta$



Since  $p$  is small, therefore  $\|\theta\|_2$  must be large to have  $p\|\theta\|_2 \geq 1$  (or  $\leq -1$ )

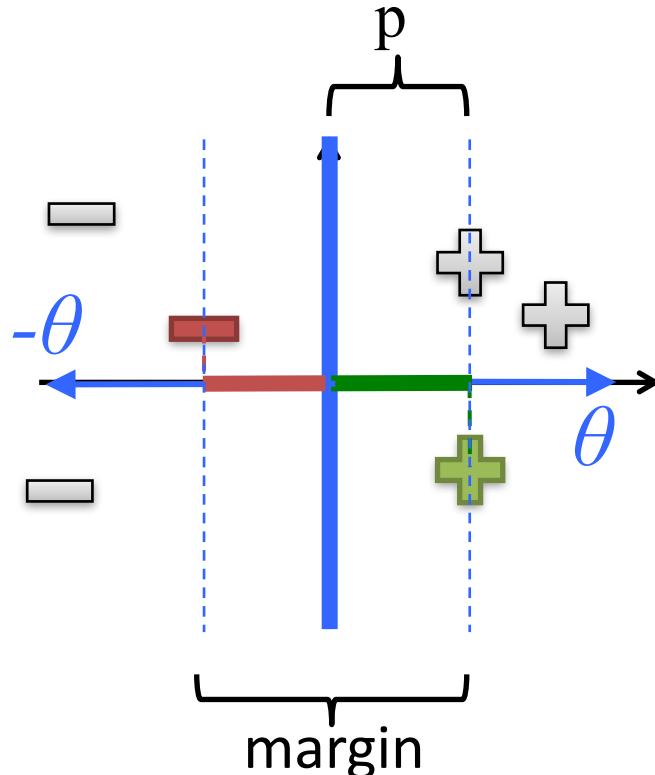


Since  $p$  is larger,  $\|\theta\|_2$  can be smaller in order to have  $p\|\theta\|_2 \geq 1$  (or  $\leq -1$ )

# Size of the Margin

For the support vectors, we have  $p\|\theta\|_2 = \pm 1$

- $p$  is the length of the projection of the SVs onto  $\theta$



Therefore,

$$p = \frac{1}{\|\theta\|_2}$$

$$\text{margin} = 2p = \frac{2}{\|\theta\|_2}$$

# The SVM Dual Problem

The primal SVM problem was given as

$$\begin{aligned} \min_{\theta} \quad & \frac{1}{2} \sum_{j=1}^d \theta_j^2 \\ \text{s.t.} \quad & y_i(\boldsymbol{\theta}^\top \mathbf{x}_i) \geq 1 \quad \forall i \end{aligned}$$

Can solve it more efficiently by taking the Lagrangian dual

- Duality is a common idea in optimization
- It transforms a difficult optimization problem into a simpler one
- Key idea: introduce slack variables  $\alpha_i$  for each constraint
  - $\alpha_i$  indicates how important a particular constraint is to the solution

# The SVM Dual Problem

- The Lagrangian is given by

$$L(\boldsymbol{\theta}, \boldsymbol{\alpha}) = \frac{1}{2} \sum_{j=1}^d \theta_j^2 - \sum_{i=1}^n \alpha_i (y_i \boldsymbol{\theta}^\top \mathbf{x} - 1)$$

s.t.  $\alpha_i \geq 0 \quad \forall i$

- We must minimize over  $\boldsymbol{\theta}$  and maximize over  $\boldsymbol{\alpha}$
- At optimal solution, partials w.r.t  $\boldsymbol{\theta}$ 's are 0

Solve by a bunch of algebra and calculus ...  
and we obtain ...

# SVM Dual Representation

$$\text{Maximize } J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

$$\sum_i \alpha_i y_i = 0$$

The decision function is given by

$$h(\mathbf{x}) = \text{sign} \left( \sum_{i \in \mathcal{SV}} \alpha_i y_i \langle \mathbf{x}, \mathbf{x}_i \rangle + b \right)$$

$$\text{where } b = \frac{1}{|\mathcal{SV}|} \sum_{i \in \mathcal{SV}} \left( y_i - \sum_{j \in \mathcal{SV}} \alpha_j y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

# Understanding the Dual

$$\text{Maximize } J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

$$\sum_i \alpha_i y_i = 0$$

Balances between the weight of constraints for different classes

Constraint weights ( $\alpha_i$ 's) cannot be negative

# Understanding the Dual

$$\text{Maximize } J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

s.t.  $\alpha_i \geq 0 \quad \forall i$

$$\sum \alpha_i y_i \leq 0$$

Points with different labels increase the sum

Points with same label decrease the sum

Measures the similarity between points

Intuitively, we should be more careful around points near the margin

# Understanding the Dual

$$\text{Maximize } J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

$$\sum_i \alpha_i y_i = 0$$

In the solution, either:

- $\alpha_i > 0$  and the constraint is tight ( $y_i(\boldsymbol{\theta}^\top \mathbf{x}_i) = 1$ )
  - point is a support vector
- $\alpha_i = 0$ 
  - point is not a support vector

# Employing the Solution

- Given the optimal solution  $\alpha^*$ , optimal weights are

$$\theta^* = \sum_{i \in SVs} \alpha_i^* y_i \mathbf{x}_i$$

- In this formulation, have *not* added  $\mathbf{x}_0 = 1$
- Therefore, we can solve one of the SV constraints

$$y_i (\theta^* \cdot \mathbf{x}_i + \theta_0) = 1$$

to obtain  $\theta_0$

- Or, more commonly, take the average solution over all support vectors

# What if Data Are Not Linearly Separable?

- Cannot find  $\theta$  that satisfies  $y_i(\theta^\top \mathbf{x}_i) \geq 1 \quad \forall i$
- Introduce slack variables  $\xi_i$ 
$$y_i(\theta^\top \mathbf{x}_i) \geq 1 - \xi_i \quad \forall i$$
- New problem:
$$\min_{\theta} \frac{1}{2} \sum_{j=1}^d \theta_j^2 + C \sum_i \xi_i$$
$$\text{s.t. } y_i(\theta^\top \mathbf{x}_i) \geq 1 - \xi_i \quad \forall i$$

# Strengths of SVMs

- Good generalization in theory
- Good generalization in practice
- Work well with few training instances
- Find globally best model
- Efficient algorithms
- Amenable to the kernel trick ...

# What if Surface is Non-Linear?

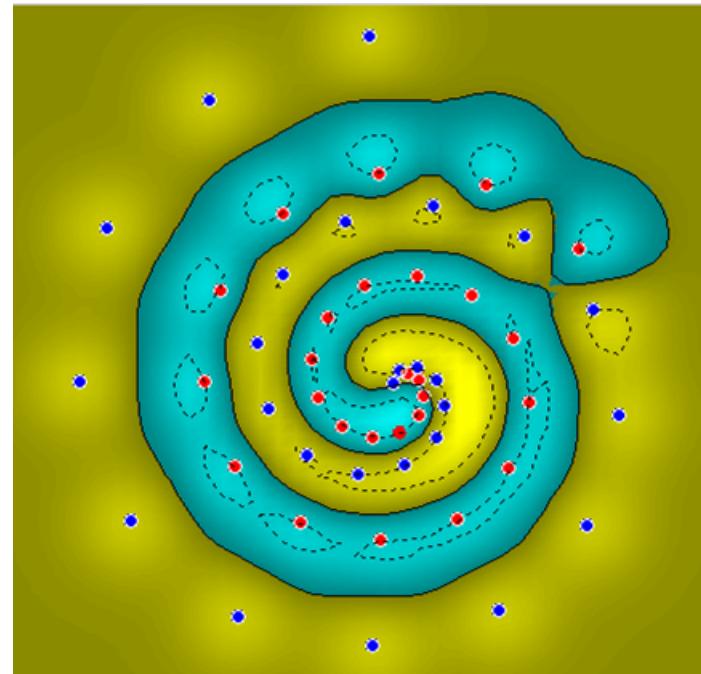
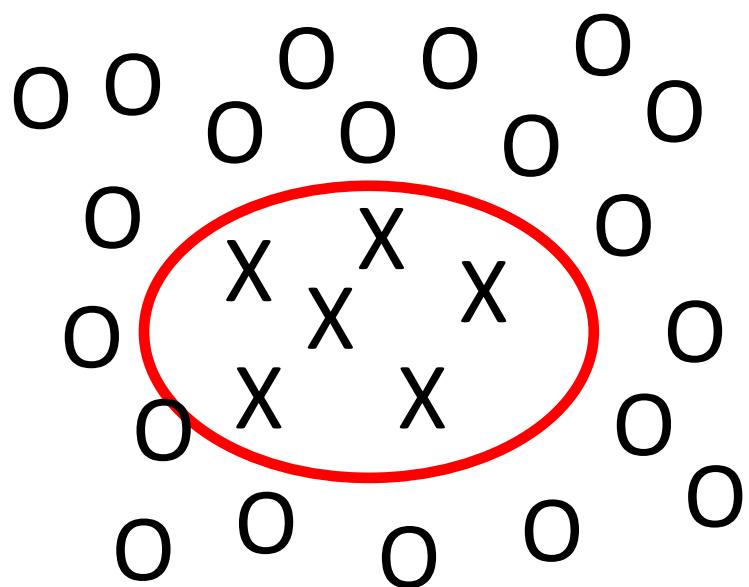
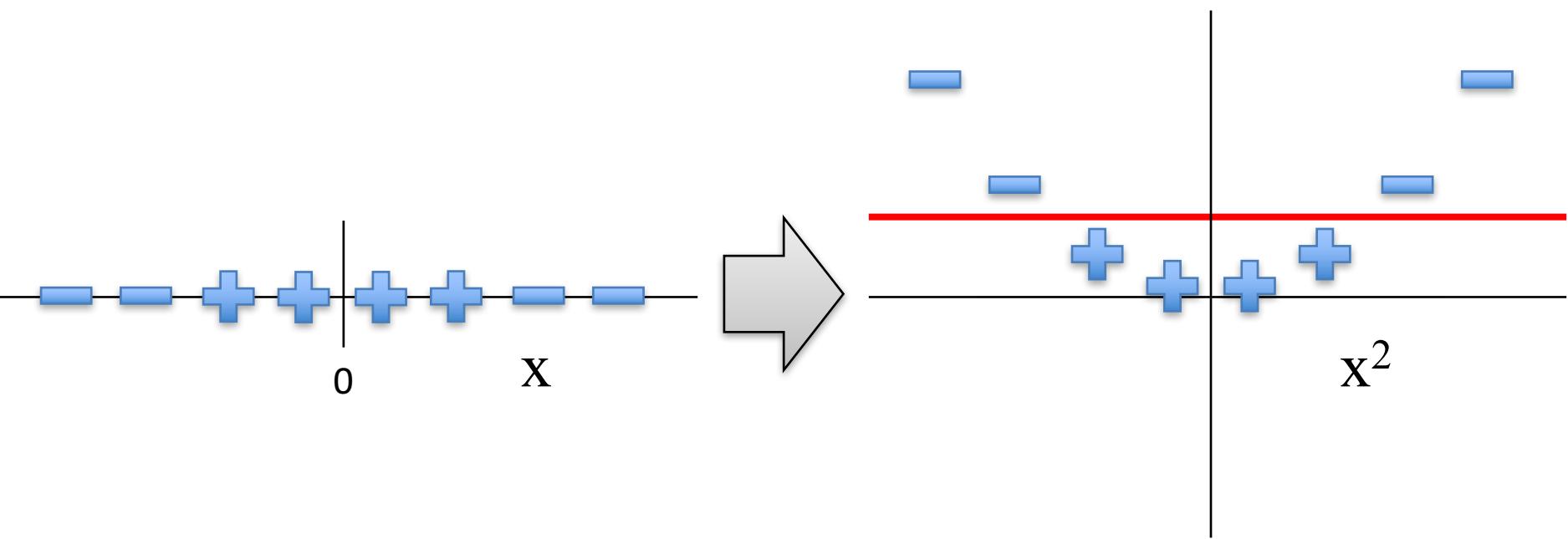


Image from <http://www.atrandomresearch.com/iclass/>

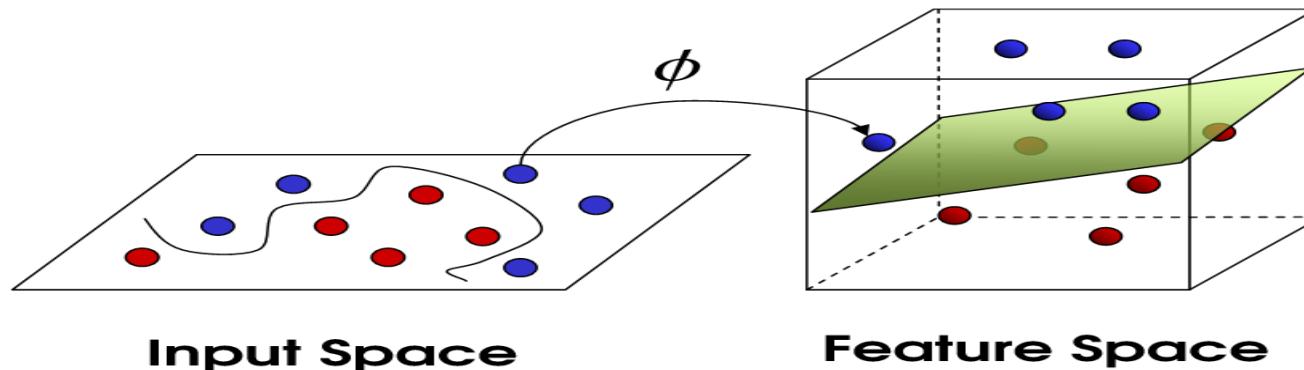
# Kernel Methods

Making the Non-Linear Linear

# When Linear Separators Fail



# Mapping into a New Feature Space



$$\Phi : \mathcal{X} \mapsto \hat{\mathcal{X}} = \Phi(\mathbf{x})$$

- For example, with  $\mathbf{x}_i \in \mathbb{R}^2$   
$$\Phi([x_{i1}, x_{i2}]) = [x_{i1}, x_{i2}, x_{i1}x_{i2}, x_{i1}^2, x_{i2}^2]$$
- Rather than run SVM on  $\mathbf{x}_i$ , run it on  $\Phi(\mathbf{x}_i)$ 
  - Find non-linear separator in input space
- What if  $\Phi(\mathbf{x}_i)$  is really big?
- Use kernels to compute it implicitly!

# Kernels

- Find kernel  $K$  such that

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$$

- Computing  $K(\mathbf{x}_i, \mathbf{x}_j)$  should be efficient, much more so than computing  $\Phi(\mathbf{x}_i)$  and  $\Phi(\mathbf{x}_j)$
- Use  $K(\mathbf{x}_i, \mathbf{x}_j)$  in SVM algorithm rather than  $\langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Remarkably, this is possible!

# The Polynomial Kernel

Let  $\mathbf{x}_i = [x_{i1}, x_{i2}]$  and  $\mathbf{x}_j = [x_{j1}, x_{j2}]$

Consider the following function:

$$\begin{aligned} K(\mathbf{x}_i, \mathbf{x}_j) &= \langle \mathbf{x}_i, \mathbf{x}_j \rangle^2 \\ &= (x_{i1}x_{j1} + x_{i2}x_{j2})^2 \\ &= (x_{i1}^2x_{j1}^2 + x_{i2}^2x_{j2}^2 + 2x_{i1}x_{i2}x_{j1}x_{j2}) \\ &= \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle \end{aligned}$$

where

$$\Phi(\mathbf{x}_i) = [x_{i1}^2, x_{i2}^2, \sqrt{2}x_{i1}x_{i2}]$$

$$\Phi(\mathbf{x}_j) = [x_{j1}^2, x_{j2}^2, \sqrt{2}x_{j1}x_{j2}]$$

# The Polynomial Kernel

- Given by  $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle^d$ 
  - $\Phi(\mathbf{x})$  contains all monomials of degree  $d$
- Useful in visual pattern recognition
  - Example:
    - 16x16 pixel image
    - $10^{10}$  monomials of degree 5
    - Never explicitly compute  $\Phi(\mathbf{x})$  !
- Variation:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^d$ 
  - Adds all lower-order monomials (degrees 1,...,d)!

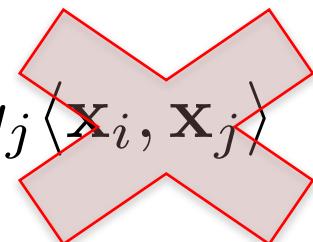
# The Kernel Trick

“Given an algorithm which is formulated in terms of a positive definite kernel  $K_1$ , one can construct an alternative algorithm by replacing  $K_1$  with another positive definite kernel  $K_2$ ”

- SVMs can use the kernel trick

# Incorporating Kernels into SVM

$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$



$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$\text{s.t. } \alpha_i \geq 0 \quad \forall i$$

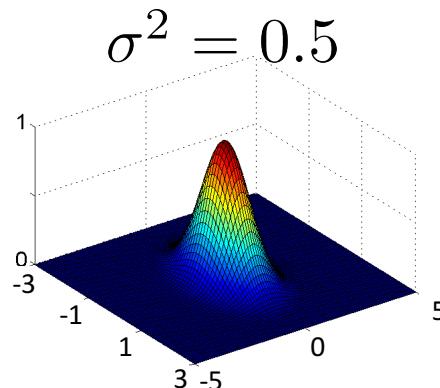
$$\sum_i \alpha_i y_i = 0$$

# The Gaussian Kernel

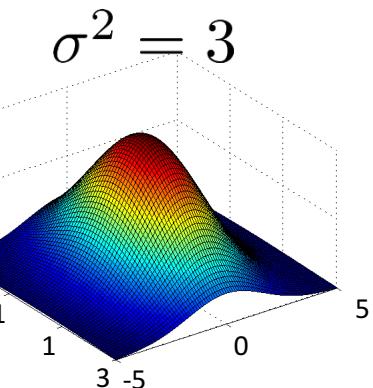
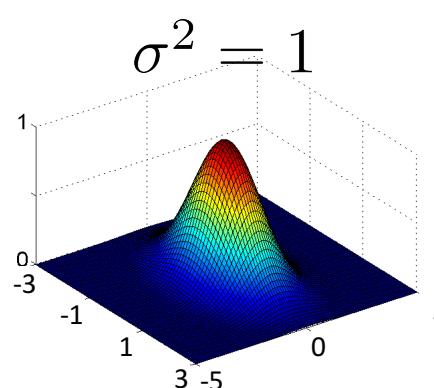
- Also called Radial Basis Function (RBF) kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

- Has value 1 when  $\mathbf{x}_i = \mathbf{x}_j$
- Value falls off to 0 with increasing distance
- Note: Need to do feature scaling before using Gaussian Kernel



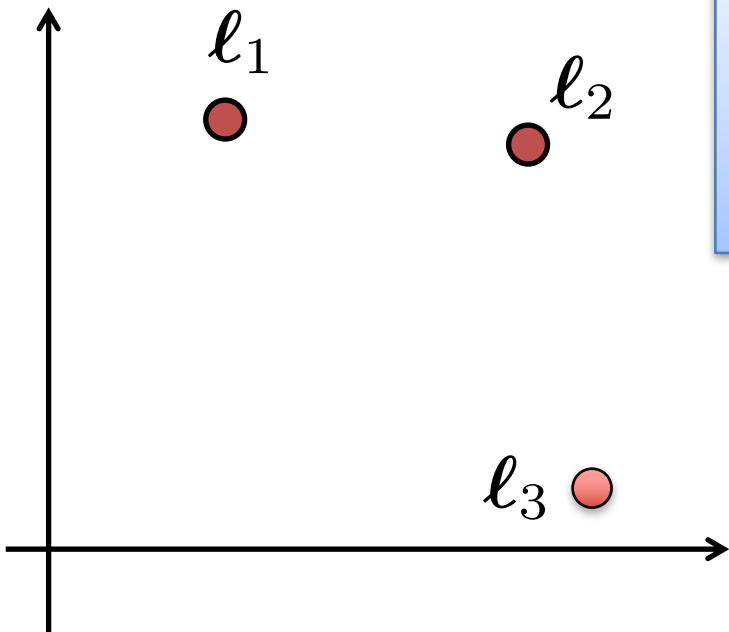
$\sigma^2 = 0.5$   
lower bias,  
higher variance



$\sigma^2 = 3$   
higher bias,  
lower variance



# Gaussian Kernel Example



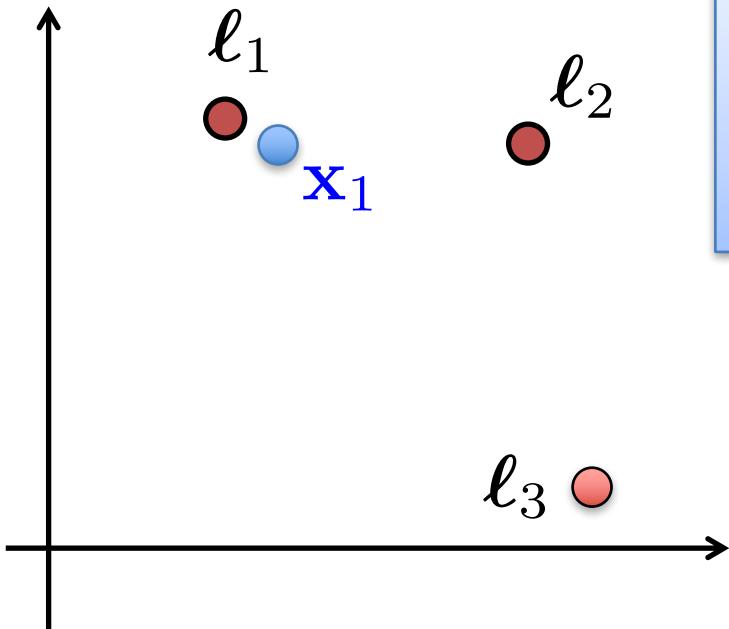
$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

$$\theta = [-0.5, 1, 1, 0]$$

Predict +1 if  $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$

# Gaussian Kernel Example



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

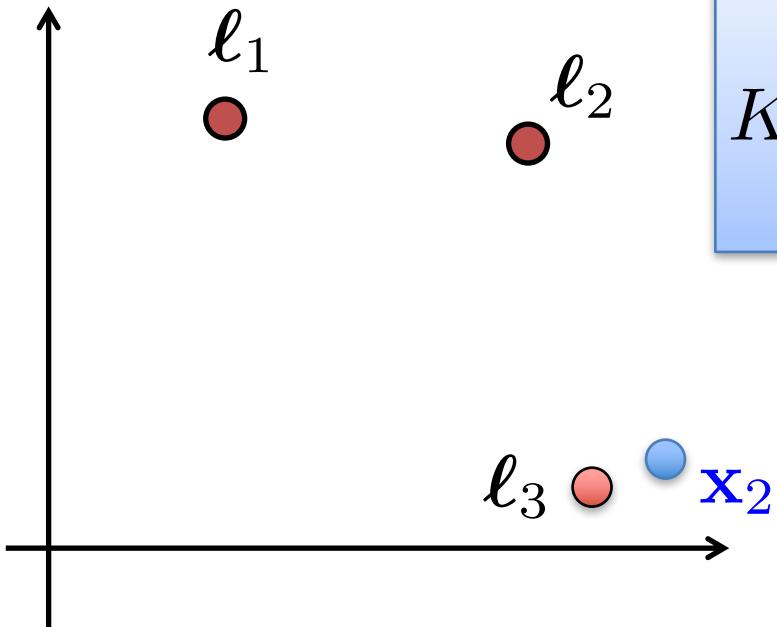
$$\theta = [-0.5, 1, 1, 0]$$

Predict +1 if  $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$

- For  $\mathbf{x}_1$ , we have  $K(\mathbf{x}_1, \ell_1) \approx 1$ , other similarities  $\approx 0$

$$\begin{aligned} \theta_0 + \theta_1(1) + \theta_2(0) + \theta_3(0) \\ = -0.5 + 1(1) + 1(0) + 0(0) \\ = 0.5 \geq 0, \text{ so predict } +1 \end{aligned}$$

# Gaussian Kernel Example



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

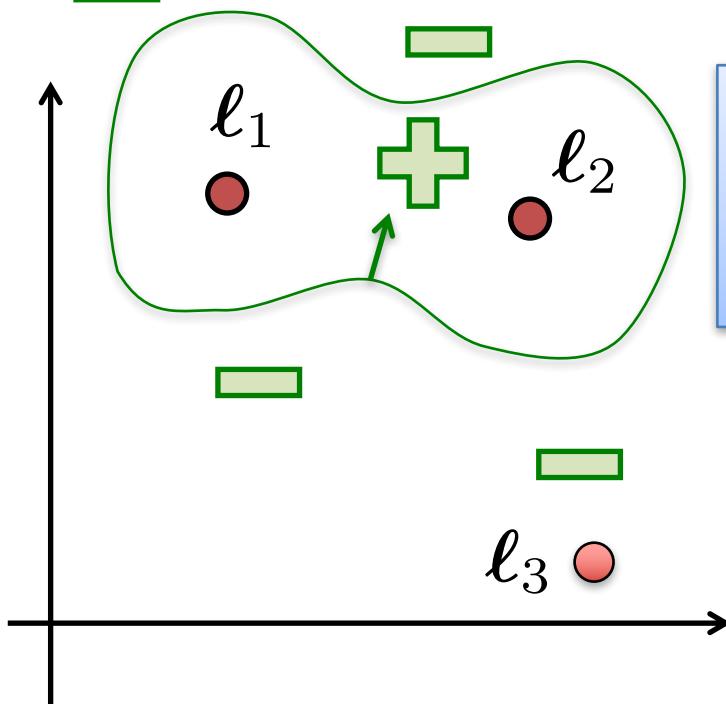
$$\theta = [-0.5, 1, 1, 0]$$

Predict +1 if  $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$

- For  $\mathbf{x}_2$ , we have  $K(\mathbf{x}_2, \ell_3) \approx 1$ , other similarities  $\approx 0$

$$\begin{aligned} \theta_0 + \theta_1(0) + \theta_2(0) + \theta_3(1) \\ = -0.5 + 1(0) + 1(0) + 0(1) \\ = -0.5 < 0, \text{ so predict -1} \end{aligned}$$

# Gaussian Kernel Example



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

$$\theta = [-0.5, 1, 1, 0]$$

Predict +1 if  $\theta_0 + \theta_1 K(\mathbf{x}, \ell_1) + \theta_2 K(\mathbf{x}, \ell_2) + \theta_3 K(\mathbf{x}, \ell_3) \geq 0$

Rough sketch of decision surface

# A Few Good Kernels...

- Linear Kernel  $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$
- Polynomial kernel  $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d$ 
  - $c \geq 0$  trades off influence of lower order terms
- Gaussian kernel  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$
- Sigmoid kernel  $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\alpha \mathbf{x}_i^\top \mathbf{x}_j + c)$

Many more...

- Cosine similarity kernel
- Chi-squared kernel
- String/tree/graph/wavelet/etc kernels

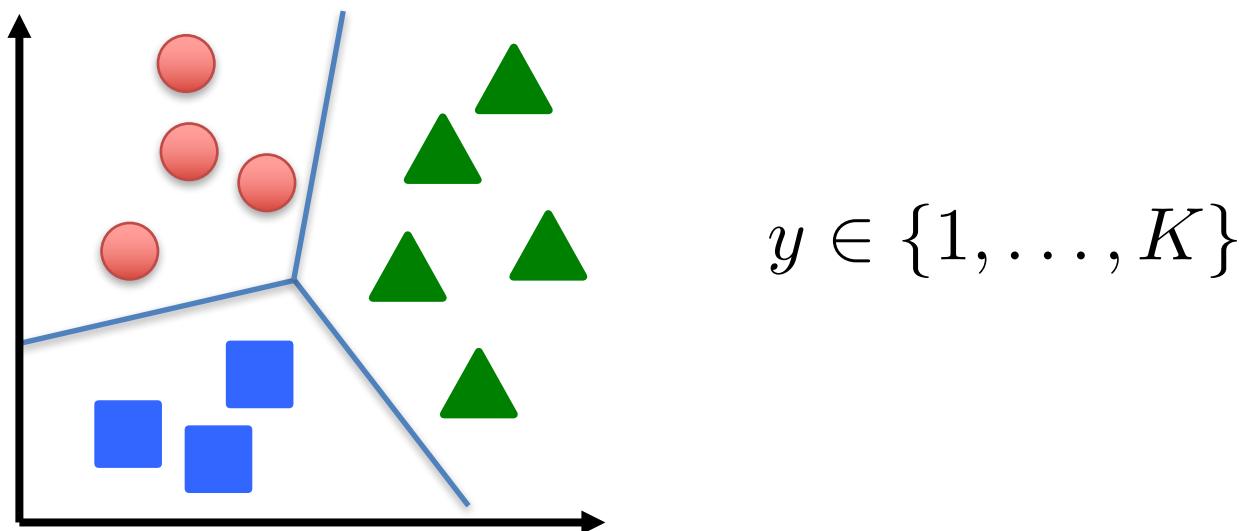
# Practical Advice for Applying SVMs

- Use SVM software package to solve for parameters
  - e.g., SVMlight, libsvm, cvx (fast!), etc.
- Need to specify:
  - Choice of parameter C
  - Choice of kernel function
    - Associated kernel parameters

e.g.,  $K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + c)^d$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

# Multi-Class Classification with SVMs



- Many SVM packages already have multi-class classification built in
- Otherwise, use one-vs-rest
  - Train  $K$  SVMs, each picks out one class from rest, yielding  $\theta^{(1)}, \dots, \theta^{(K)}$
  - Predict class  $i$  with largest  $(\theta^{(i)})^\top \mathbf{x}$

# SVMs vs Logistic Regression

## (Advice from Andrew Ng)

$n = \# \text{ training examples}$      $d = \# \text{ features}$

If  $d$  is large (relative to  $n$ ) (e.g.,  $d > n$  with  $d = 10,000$ ,  $n = 10-1,000$ )

- Use logistic regression or SVM with a linear kernel

If  $d$  is small (up to 1,000),  $n$  is intermediate (up to 10,000)

- Use SVM with Gaussian kernel

If  $d$  is small (up to 1,000),  $n$  is large (50,000+)

- Create/add more features, then use logistic regression or SVM without a kernel

Neural networks likely to work well for most of these settings, but may be slower to train

# Conclusion

- SVMs find optimal linear separator
- The kernel trick makes SVMs learn non-linear decision surfaces
- Strength of SVMs:
  - Good theoretical and empirical performance
  - Supports many types of kernels
- Disadvantages of SVMs:
  - “Slow” to train/predict for huge data sets (but relatively fast!)
  - Need to choose the kernel (and tune its parameters)