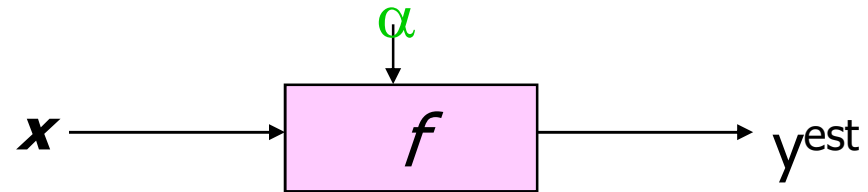


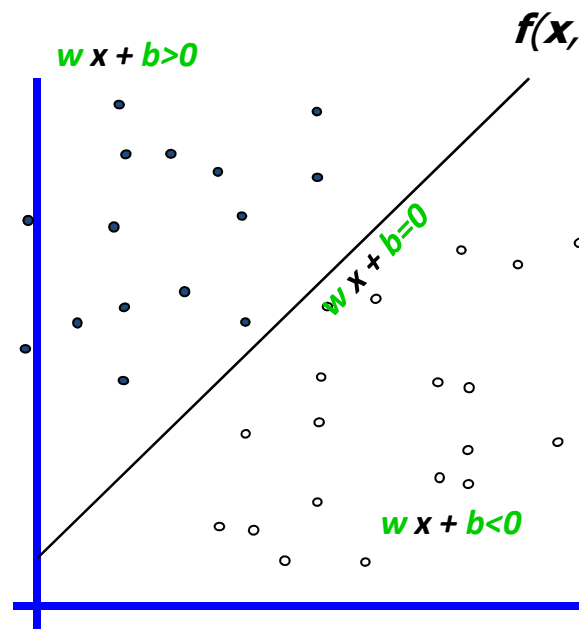
# Support Vector Machines

- **Duration:** 2 hrs
- Support Vector Machine
  - Linear classifier
  - Maximum margins
  - Soft margin classification
  - Non-linear SVM
  - Kernel functions
  - Multiclass classification

# Linear Classifiers



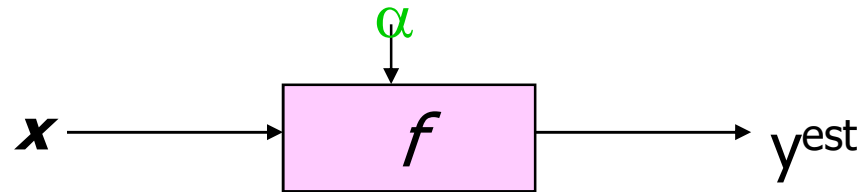
- denotes +1
- denotes -1



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w}\mathbf{x} + b)$$

How would you classify this data?

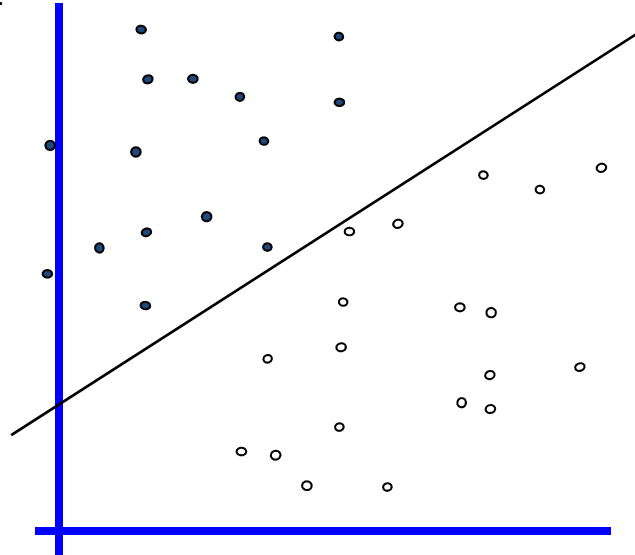
# Linear Classifiers



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

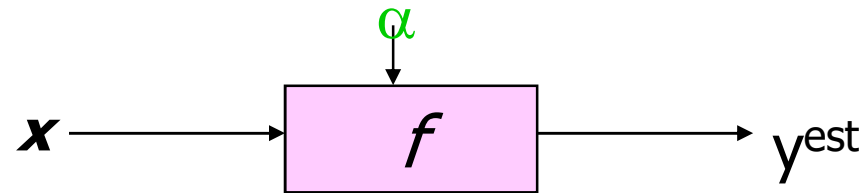
• denotes +1

° denotes -1



How would you  
classify this data?

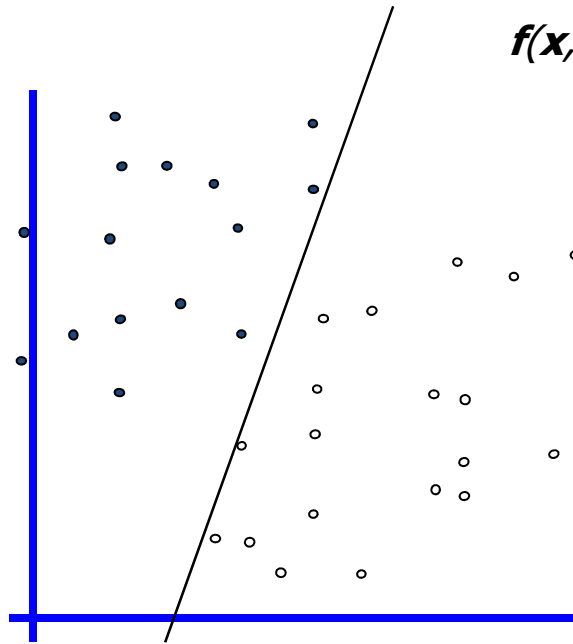
# Linear Classifiers



$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

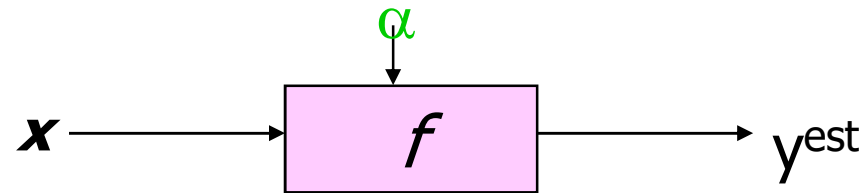
• denotes +1

◦ denotes -1

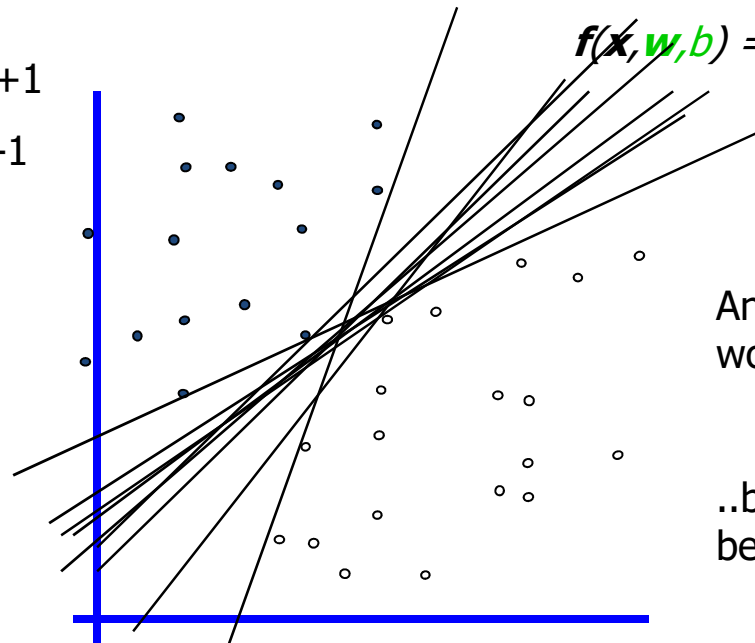


How would you  
classify this data?

# Linear Classifiers



• denotes +1  
○ denotes -1

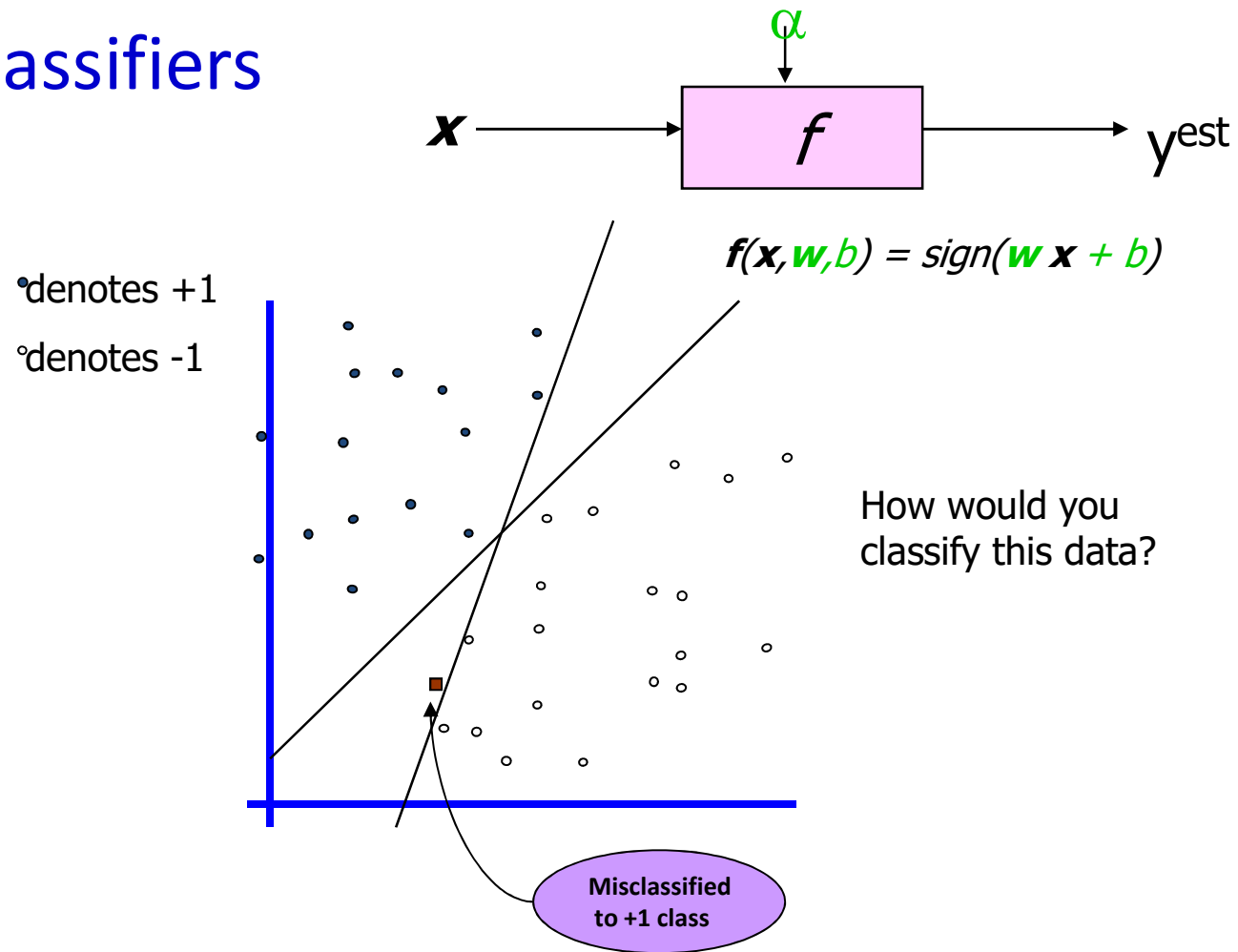


$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

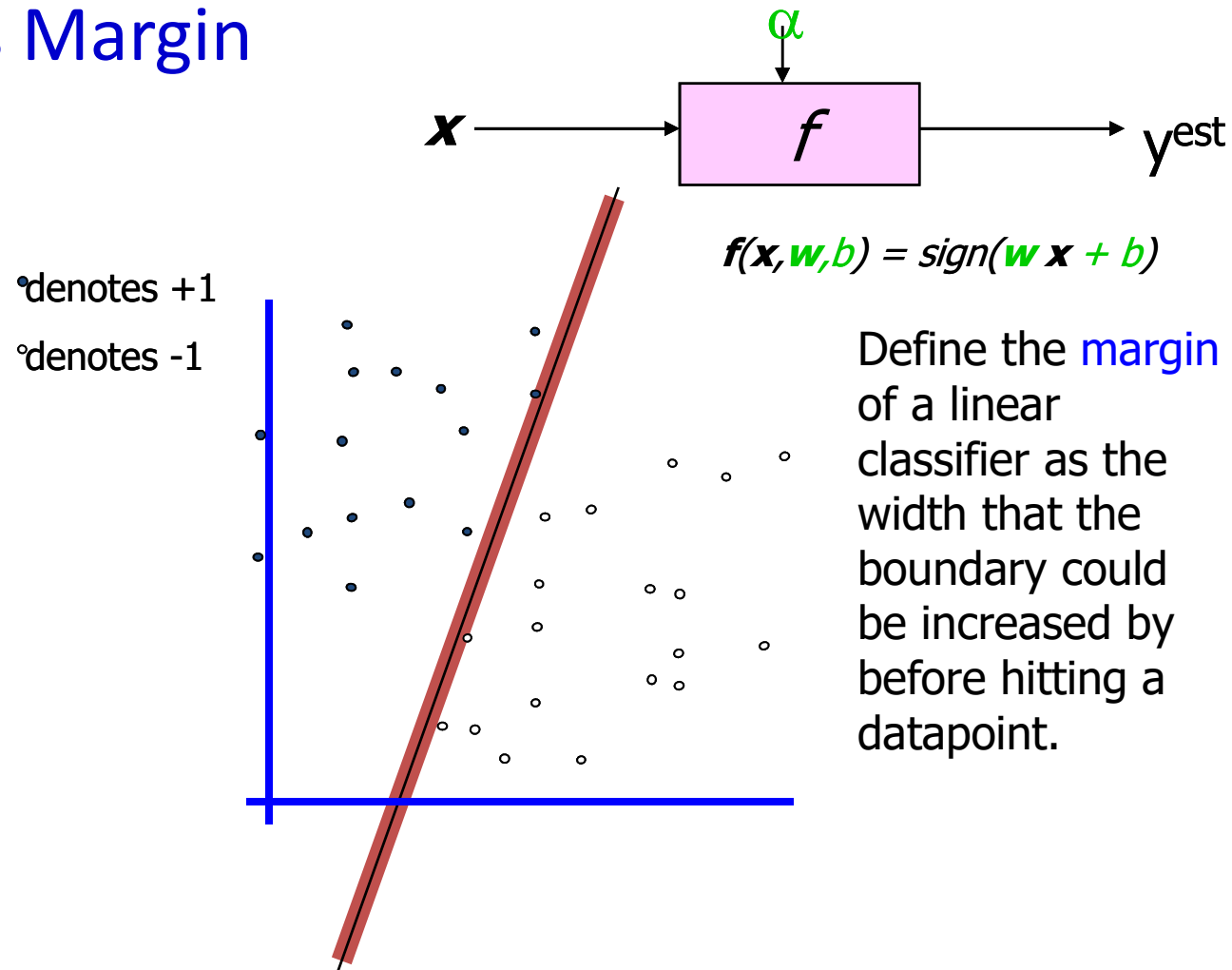
Any of these  
would be fine..

..but which is  
best?

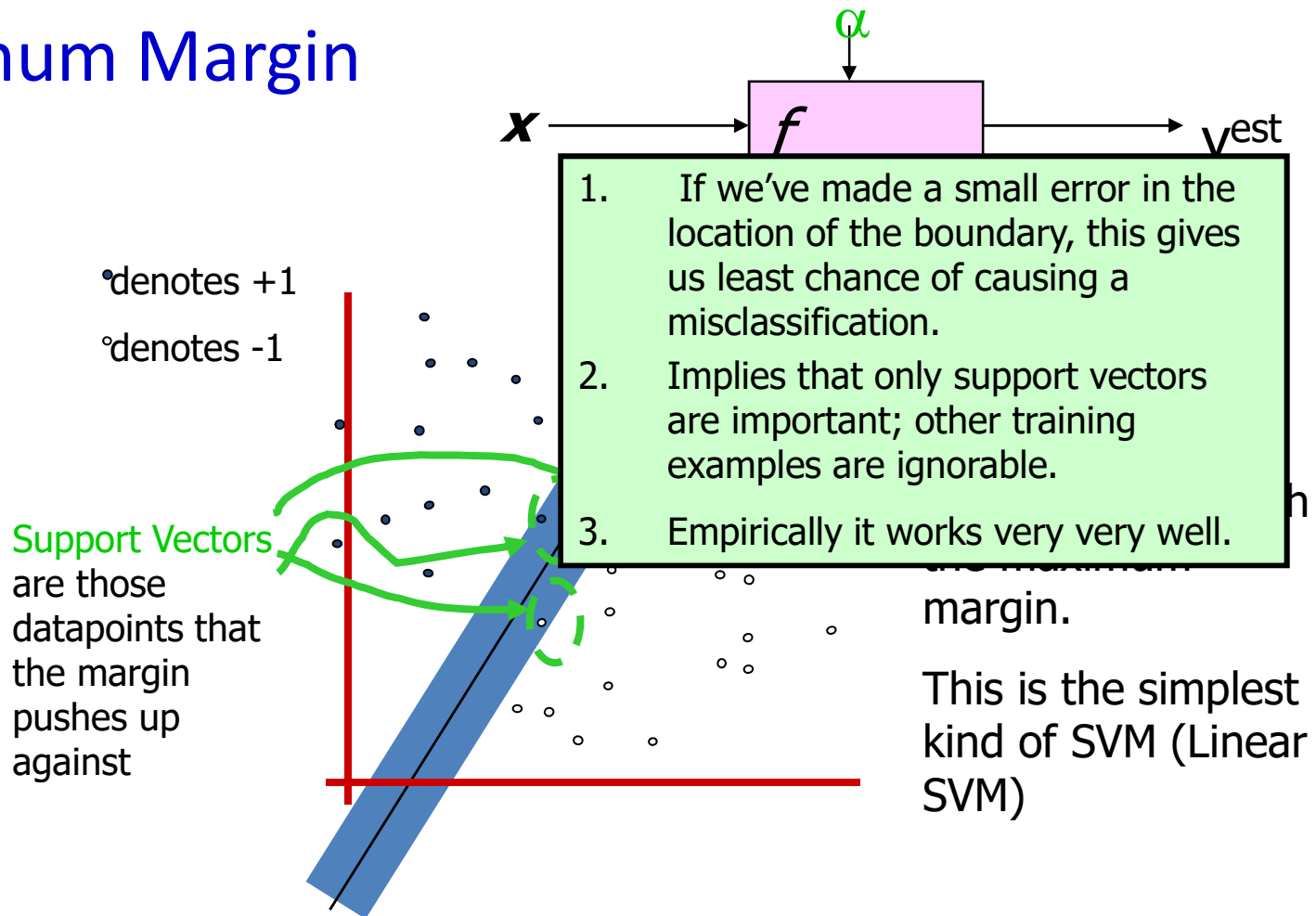
# Linear Classifiers



# Classifiers Margin

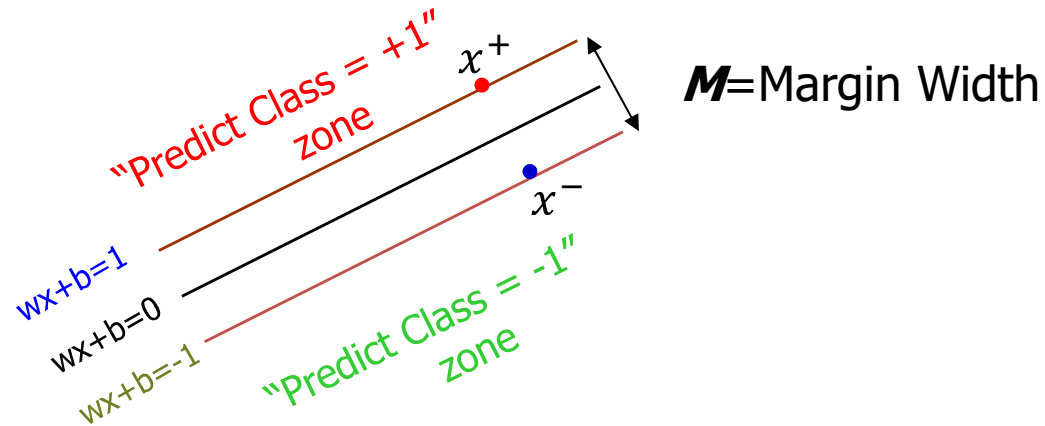


# Maximum Margin





# Linear SVM Mathematically



What we know:

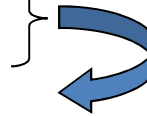
- $w \cdot x^+ + b = +1$
- $w \cdot x^- + b = -1$
- $w \cdot (x^+ - x^-) = 2$
- $x^+ - x^- = \lambda w$  for some value of

$$M = |x^+ - x^-| = \lambda |w|$$

$$= \frac{2}{w \cdot w} |w| = \frac{2}{|w|}$$

# Linear SVM Mathematically

- Goal: 1) **Correctly classify all training data**

$$\left. \begin{array}{ll} wx_i + b \geq 1 & \text{if } y_i = +1 \\ wx_i + b \leq 1 & \text{if } y_i = -1 \end{array} \right\} \text{for all } i$$


2) **Maximize the Margin**  $M = \frac{2}{|w|}$   
**same as minimize**  $\frac{1}{2} w^t w$

- We can formulate a Quadratic Optimization Problem and solve for w and b

Minimize  $\Phi(w) = \frac{1}{2} w^t w$

subject to  $y_i (wx_i + b) \geq 1 \quad \forall i$

# The Optimization Problem Solution

- The solution has the form:

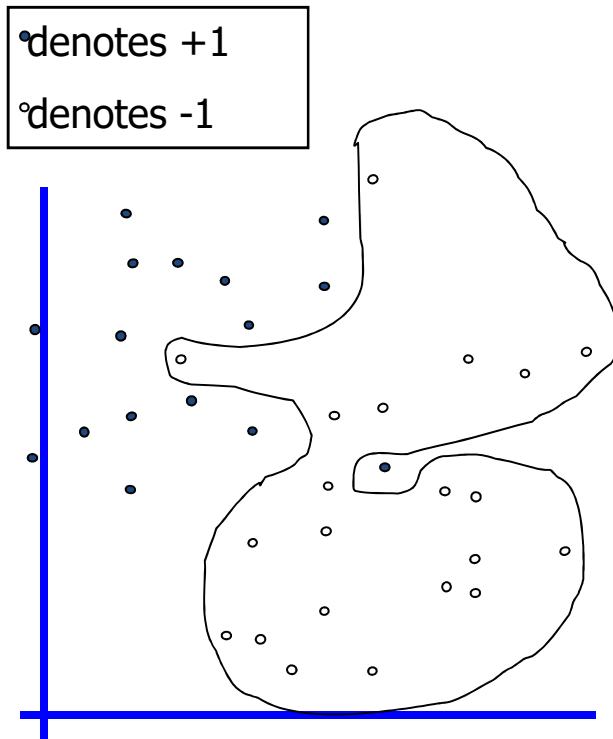
$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = y_k - \mathbf{w}^T \mathbf{x}_k \text{ for any } \mathbf{x}_k \text{ such that } \alpha_k \neq 0$$

- Each non-zero  $\alpha_i$  indicates that corresponding  $\mathbf{x}_i$  is a support vector.
- Then the classifying function will have the form:

$$f(\mathbf{x}) = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$$

- Notice that it relies on an *inner product* between the test point  $\mathbf{x}$  and the support vectors  $\mathbf{x}_i$  – we will return to this later.
- Also keep in mind that solving the optimization problem involved computing the inner products  $\mathbf{x}_i^T \mathbf{x}_j$  between all pairs of training points.

## Dataset with noise

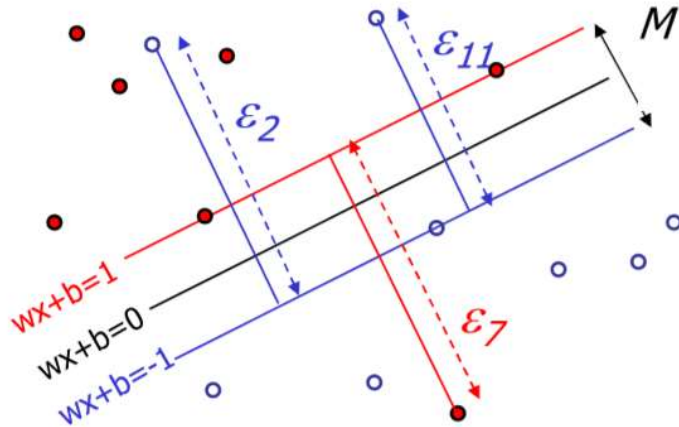


- Hard Margin: So far we require all data points be classified correctly
  - No training error
- What if the training set is noisy?
- - Solution 1: use very powerful kernels

**OVERFITTING!**

# Soft Margin Classification

Slack variables  $\xi_i$  can be added to allow misclassification of difficult or noisy examples.



What should our quadratic optimization criterion be?

Minimize

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

# Hard Margin v.s. Soft Margin

- **The old formulation:**

Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$  is minimized and for all  $\{(\mathbf{x}_i, y_i)\}$   
 $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

- **The new formulation incorporating slack variables:**

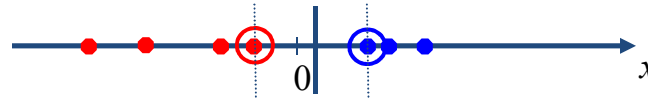
Find  $\mathbf{w}$  and  $b$  such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i$  is minimized and for all  $\{(\mathbf{x}_i, y_i)\}$   
 $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i$  and  $\xi_i \geq 0$  for all  $i$

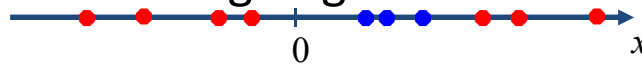
- **Parameter  $C$  can be viewed as a way to control overfitting.**

# Non-linear SVMs

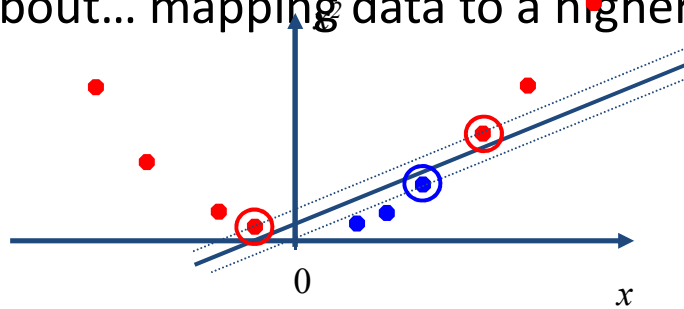
- Datasets that are linearly separable with some noise work out great:



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

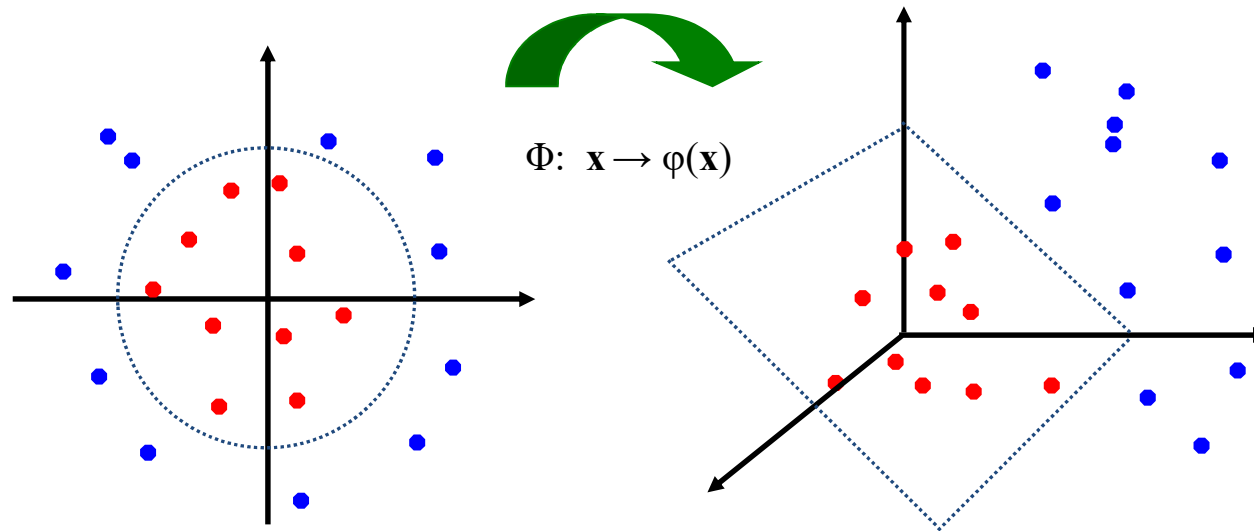


- How about... mapping data to a higher-dimensional space:



## 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:





## The “Kernel Trick”

- The linear classifier relies on dot product between vectors  $K(x_i, x_j) = x_i^T x_j$
- 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.
- Example:

2-dimensional vectors  $x = [x_1 \ x_2]$ ; let  $K(x_i, x_j) = (1 + x_i^T x_j)^2$ ,

Need to show that  $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$ :

$$\begin{aligned}
 K(x_i, x_j) &= (1 + x_i^T x_j)^2, \\
 &= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \\
 &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\
 &= \phi(x_i)^T \phi(x_j), \quad \text{where } \phi(x) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2]
 \end{aligned}$$

## Examples of Kernel Functions

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial of power  $p$ :  $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
- Gaussian (radial-basis function network):  
$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$
- Sigmoid:  $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^T \mathbf{x}_j + \beta_1)$

Use SVM software package (e.g. liblinear, libsvm, ...)

$\theta$

Need to specify:

Choice of parameter C.

Choice of kernel (similarity function):

E.g. No kernel (“linear kernel”)

Gaussian kernel:

$$f_i = \exp \left( -\frac{\|x - l^{(i)}\|^2}{2\sigma^2} \right), \text{ where } l^{(i)} = x^{(i)}.$$

Need to choose  $\sigma^2$ .

## Kernel (similarity) functions:

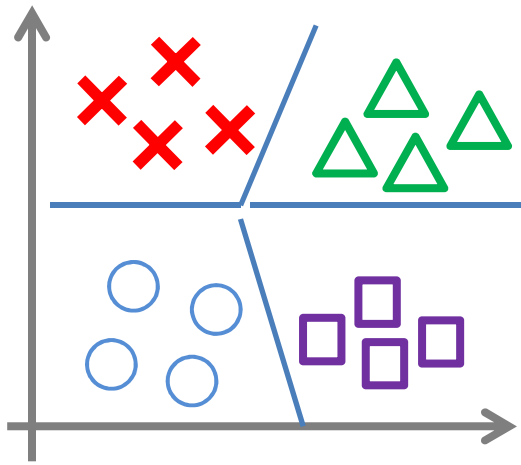
```
function f = kernel(x1,x2)
```

$$f = \exp \left( -\frac{\| \mathbf{x1} - \mathbf{x2} \|^2}{2\sigma^2} \right)$$

```
return
```

Note: Do perform feature scaling before using the Gaussian kernel.

## Multi-class classification



$$y \in \{1, 2, 3, \dots, K\}$$

Many SVM packages already have built-in multi-class classification functionality.

Otherwise, use one-vs.-all method. (Train  $K$  SVMs, one to distinguish  $y = i$  from the rest, for  $i = 1, 2, \dots, K$ )

## Logistic regression vs. SVMs

$n$  = number of features ( $x \in \mathbb{R}^{n+1}$ ),  $m$  = number of training examples

If  $n$  is large (relative to  $m$ ): (E.g,  $n \geq m, n = 10\,000, m = 10 - 1000$ )

Use logistic regression, or SVM without a kernel (“linear kernel”)

If  $n$  is small,  $m$  is intermediate: ( $n = 1 - 1000, m = 10 - 10000$ )

Use SVM with Gaussian kernel

If  $n$  is small,  $m$  is large: ( $n = 1 - 1000, m = 50000$ )

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

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

# References

- Machine Learning, Andrew Ng, coursera.org
- SVM tutorial, Prof. Andrew Moore <http://www.cs.cmu.edu/~awm/tutorials>