

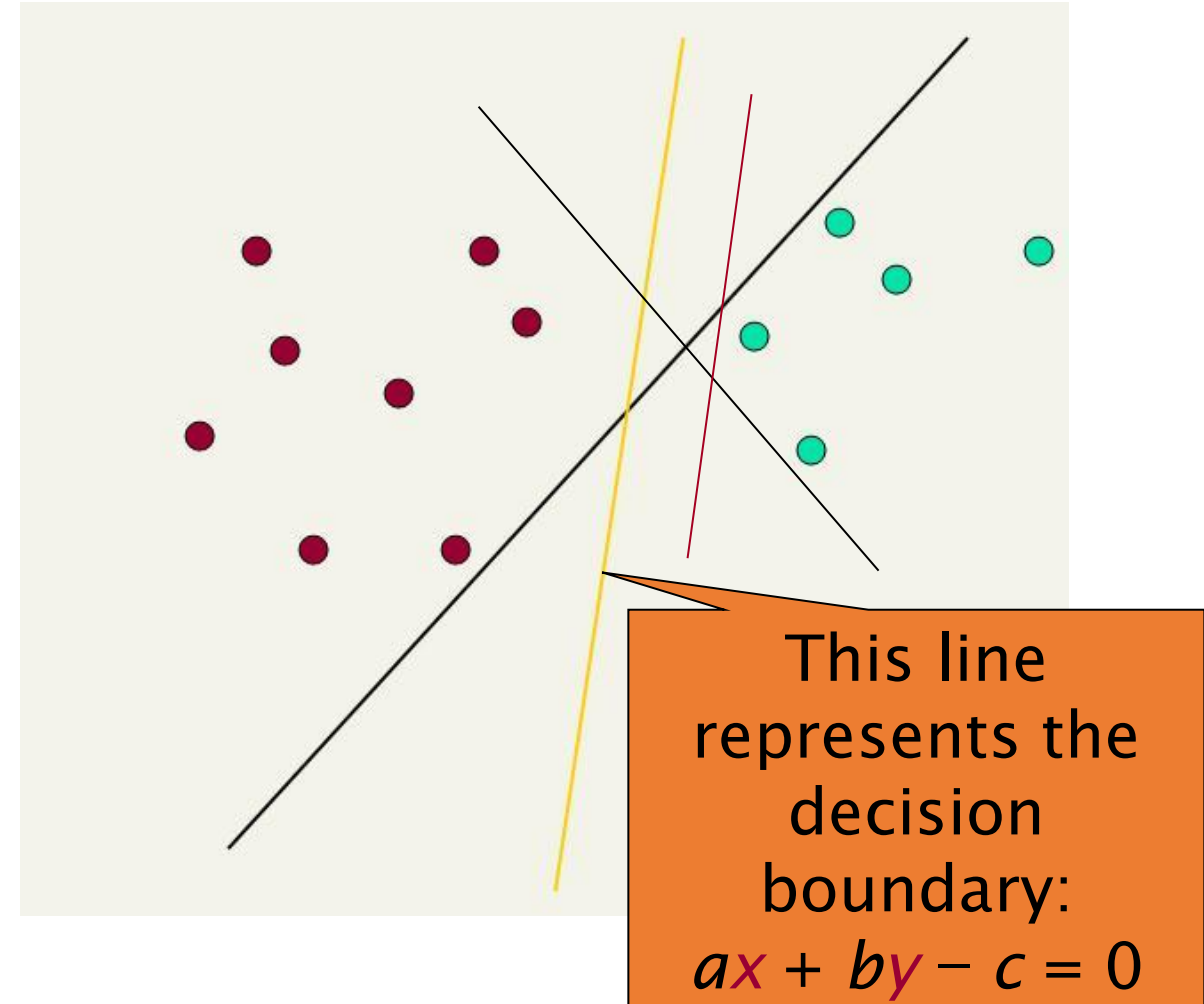
# Machine Learning

## Support Vector Machines

Edwin Puertas, Ph.D(c).  
[epuerta@utb.edu.co](mailto:epuerta@utb.edu.co)

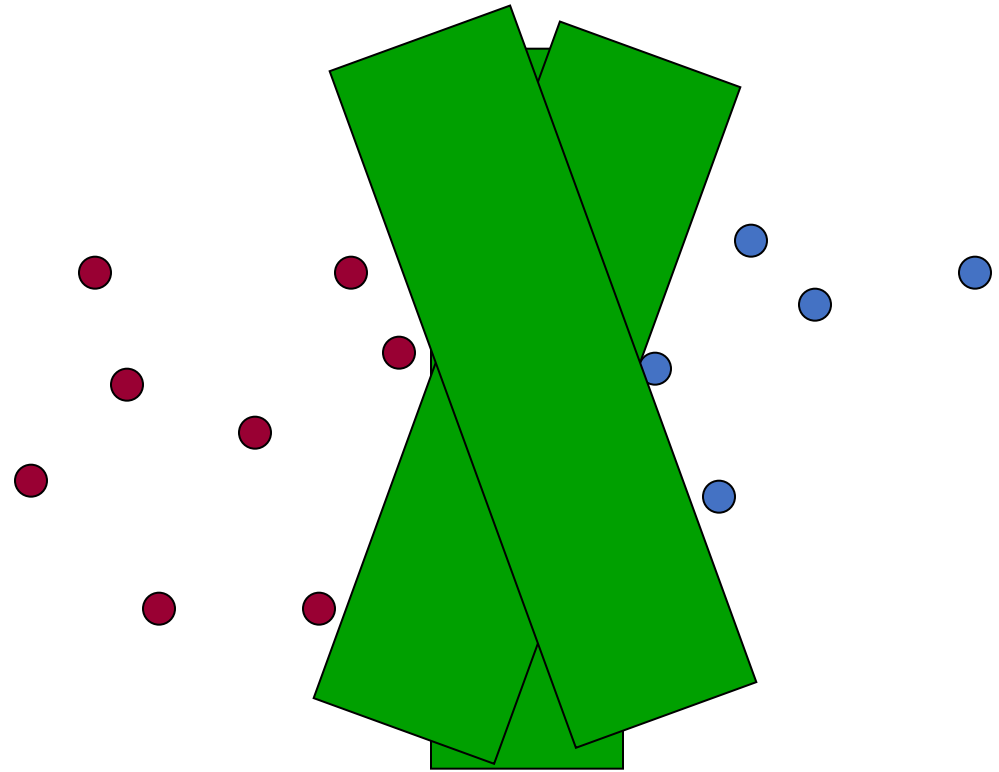
# Linear classifiers: Which Hyperplane?

- Lots of possible solutions for  $a$ ,  $b$ ,  $c$ .
- Some methods find a separating hyperplane, but not the optimal one [according to some criterion of expected goodness]
  - E.g., perceptron
- Support Vector Machine (SVM) finds an optimal\* solution.
  - Maximizes the distance between the hyperplane and the “difficult points” close to decision boundary
  - One intuition: if there are no points near the decision surface, then there are no very uncertain classification decisions



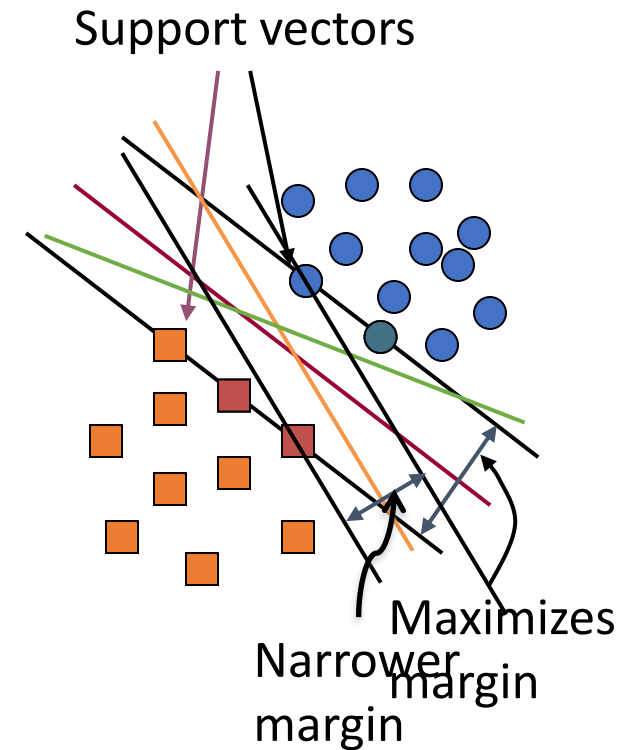
# Another intuition

If you have to place a fat separator between classes, you have less choices, and so the capacity of the model has been decreased



# Support Vector Machine (SVM)

- SVMs maximize the *margin* around the separating hyperplane.
  - A.k.a. large margin classifiers
- The decision function is fully specified by a subset of training samples, *the support vectors*.
- Solving SVMs is a *quadratic programming* problem
- Seen by many as the most successful current text classification method\*



\*but other discriminative methods often perform very similarly

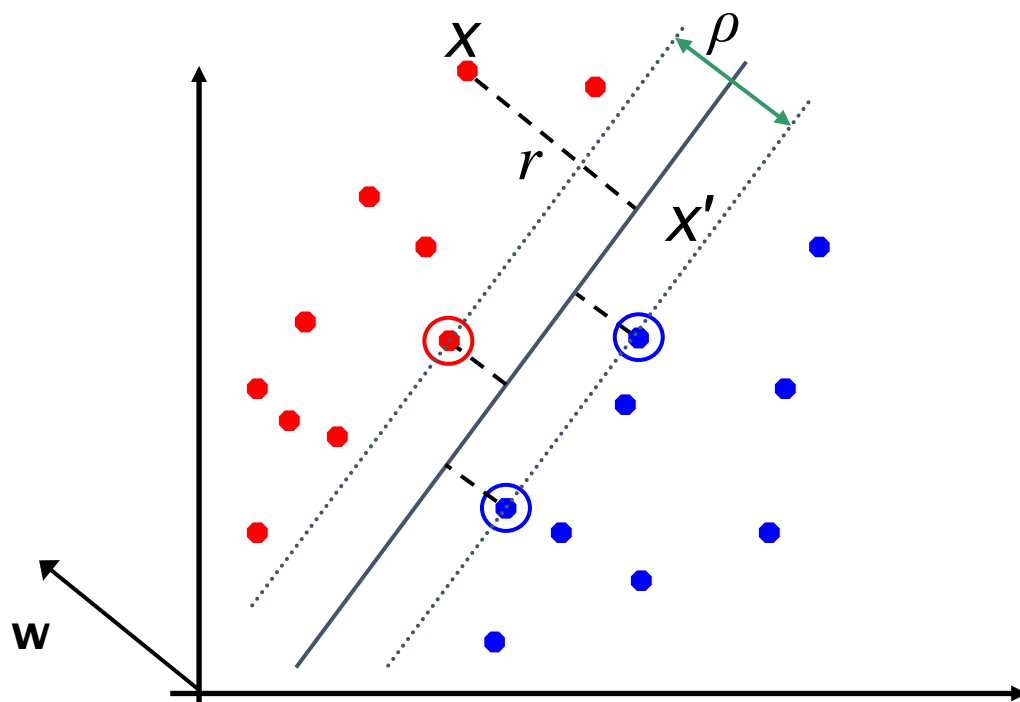
# Maximum Margin: Formalization

- $\mathbf{w}$ : decision hyperplane normal vector
- $\mathbf{x}_i$ : data point  $i$
- $y_i$ : class of data point  $i$  (+1 or -1)    NB: Not 1/0
- Classifier is:  $f(\mathbf{x}_i) = \text{sign}(\mathbf{w}^T \mathbf{x}_i + b)$
- Functional margin of  $\mathbf{x}_i$  is:  $y_i (\mathbf{w}^T \mathbf{x}_i + b)$ 
  - But note that we can increase this margin simply by scaling  $\mathbf{w}$ ,  $b$ ....
- Functional margin of dataset is twice the minimum functional margin for any point
  - The factor of 2 comes from measuring the whole width of the margin

# Geometric Margin

- Distance from example to the separator is
- Examples closest to the hyperplane are **support vectors**.
- **Margin**  $\rho$  of the separator is the width of separation between support vectors of classes.

$$r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$$



Derivation of finding  $r$ :  
 Dotted line  $\mathbf{x}' - \mathbf{x}$  is perpendicular to decision boundary so parallel to  $\mathbf{w}$ .  
 Unit vector is  $\mathbf{w}/|\mathbf{w}|$ , so line is  $r\mathbf{w}/|\mathbf{w}|$ .  
 $\mathbf{x}' = \mathbf{x} - yr\mathbf{w}/|\mathbf{w}|$ .  
 $\mathbf{x}'$  satisfies  $\mathbf{w}^T \mathbf{x}' + b = 0$ .  
 So  $\mathbf{w}^T (\mathbf{x} - yr\mathbf{w}/|\mathbf{w}|) + b = 0$   
 Recall that  $|\mathbf{w}| = \sqrt{\mathbf{w}^T \mathbf{w}}$ .  
 So  $\mathbf{w}^T \mathbf{x} - yr|\mathbf{w}| + b = 0$   
 So, solving for  $r$  gives:  
 $r = y(\mathbf{w}^T \mathbf{x} + b)/|\mathbf{w}|$

# Linear SVM Mathematically

The linearly separable case

- Assume that all data is at least distance 1 from the hyperplane, then the following two constraints follow for a training set  $\{(\mathbf{x}_i, y_i)\}$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \quad \text{if } y_i = 1$$

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{if } y_i = -1$$

- For support vectors, the inequality becomes an equality
- Then, since each example's distance from the hyperplane is

$$r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$$

- The margin is:  $r = \frac{2}{\|\mathbf{w}\|}$

# Linear Support Vector Machine (SVM)

- **Hyperplane**

- $\mathbf{w}^T \mathbf{x} + b = 0$

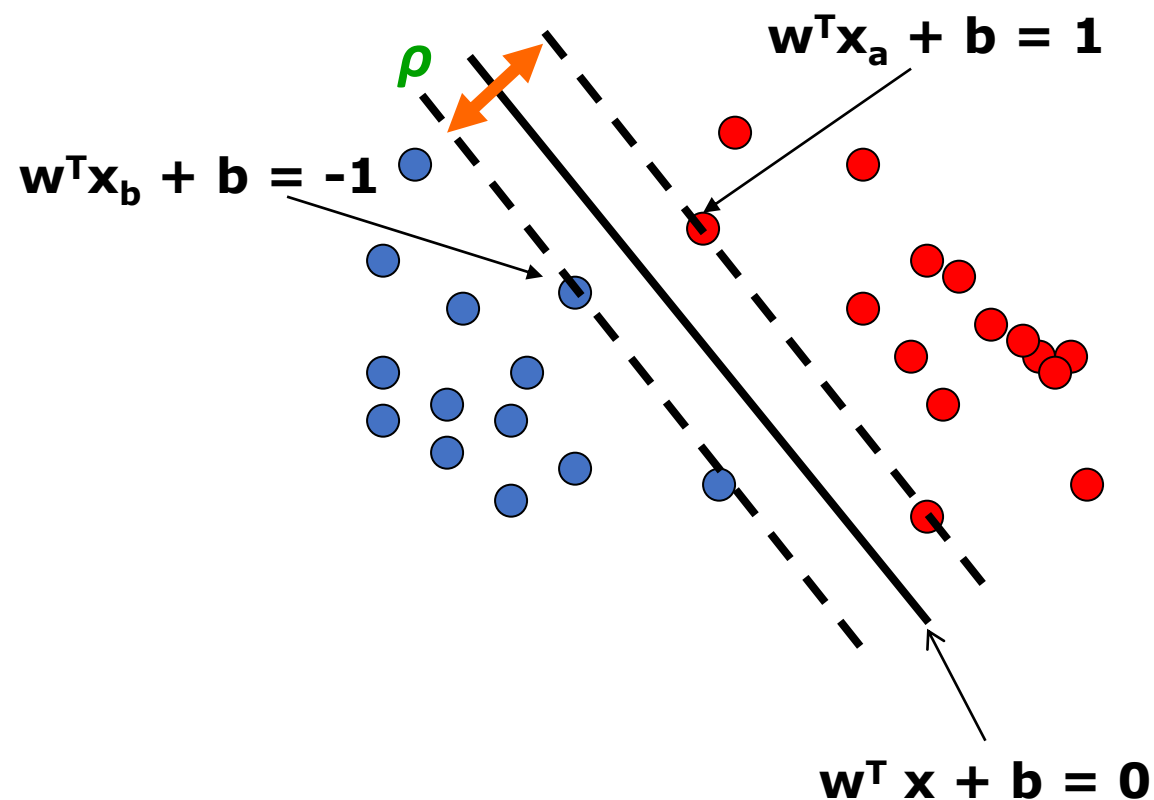
- **Extra scale constraint:**

- $\min_{i=1,\dots,n} |\mathbf{w}^T \mathbf{x}_i + b| = 1$

- This implies:

- $\mathbf{w}^T (\mathbf{x}_a - \mathbf{x}_b) = 2$

- $\rho = \|\mathbf{x}_a - \mathbf{x}_b\|_2 = 2 / \|\mathbf{w}\|_2$





# Linear SVMs Mathematically (cont.)

- Then we can formulate the *quadratic optimization problem*:

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

$$r = \frac{2}{\|\mathbf{w}\|} \quad \text{is maximized; and for all } \{(\mathbf{x}_i, y_i)\}$$

$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \text{ if } y_i=1; \quad \mathbf{w}^T \mathbf{x}_i + b \leq -1 \text{ if } y_i=-1$$

- A better formulation ( $\min \|\mathbf{w}\| = \max 1/\|\mathbf{w}\|$ ):

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

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

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

# Solving the Optimization Problem

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

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

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

- This is now optimizing a *quadratic* function subject to *linear* constraints
- Quadratic optimization problems are a well-known class of mathematical programming problem, and many (intricate) algorithms exist for solving them (with many special ones built for SVMs)
- The solution involves constructing a *dual problem* where a *Lagrange multiplier*  $\alpha_i$  is associated with every constraint in the primary problem:

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

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

(2)  $\alpha_i \geq 0$  for all  $\alpha_i$

# The Optimization Problem Solution

- The solution has the form:

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

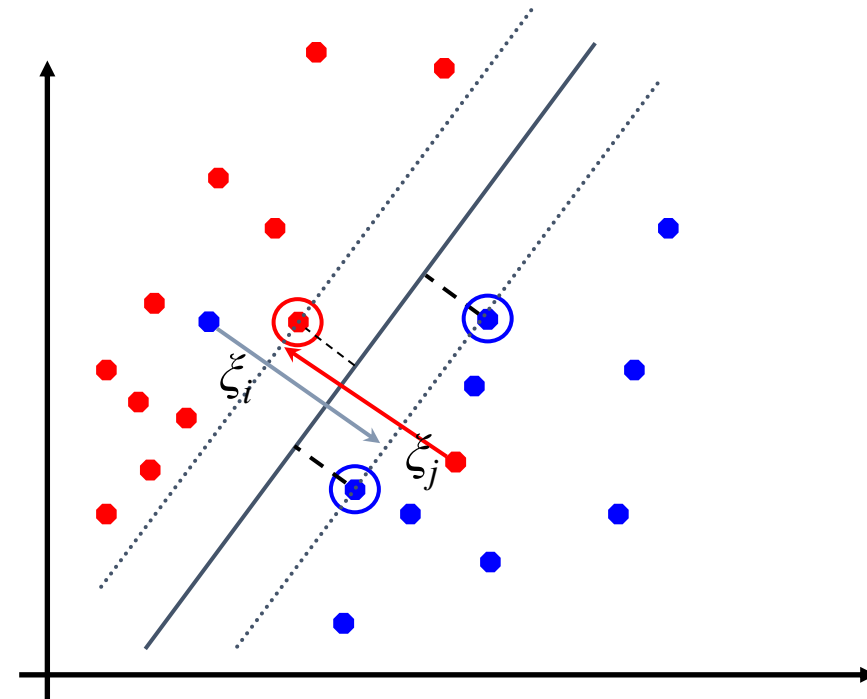
- Each non-zero  $\alpha_i$  indicates that corresponding  $\mathbf{x}_i$  is a support vector.
- Then the classifying function will have the form:

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

- 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.

# Soft Margin Classification

- If the training data is not linearly separable, *slack variables*  $\xi_i$  can be added to allow misclassification of difficult or noisy examples.
- Allow some errors
  - Let some points be moved to where they belong, at a cost
- Still, try to minimize training set errors, and to place hyperplane “far” from each class (large margin)



# Soft Margin Classification Mathematically

- 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
  - A regularization term

# Soft Margin Classification – Solution

- The dual problem for soft margin classification:

Find  $\alpha_1 \dots \alpha_N$  such that

$Q(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

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

(2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$

- Neither slack variables  $\xi_i$  nor their Lagrange multipliers appear in the dual problem!
- Again,  $\mathbf{x}_i$  with non-zero  $\alpha_i$  will be support vectors.
- Solution to the dual problem is:

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i$$

$$b = y_k (1 - \xi_k) - \mathbf{w}^T \mathbf{x}_k \text{ where } k = \underset{k'}{\operatorname{argmax}} \alpha_k$$

$\mathbf{w}$  is not needed explicitly for classification!

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

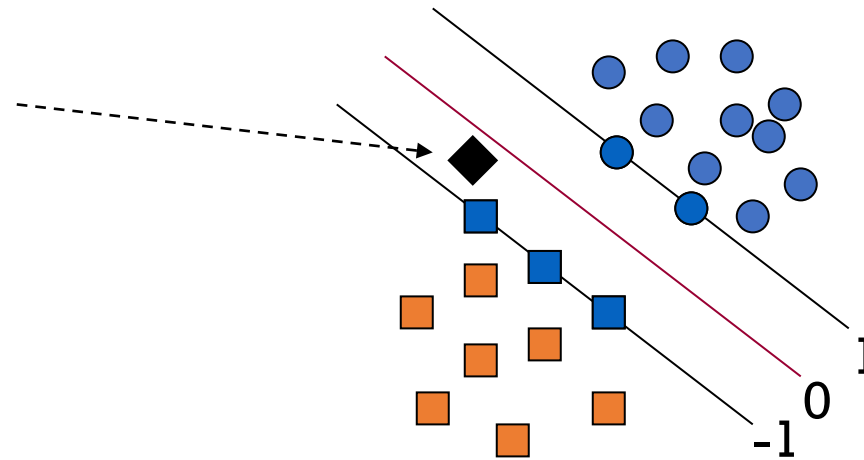
# Classification with SVMs

- Given a new point  $\mathbf{x}$ , we can score its projection onto the hyperplane normal:
  - I.e., compute score:  $\mathbf{w}^T \mathbf{x} + b = \sum \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b$ 
    - Decide class based on whether  $<$  or  $>$  0
- Can set confidence threshold  $t$ .

Score  $> t$ : yes

Score  $< -t$ : no

Else: don't know



# Linear SVMs: Summary

- The classifier is a *separating hyperplane*.
- The most “important” training points are the support vectors; they define the hyperplane.
- Quadratic optimization algorithms can identify which training points  $\mathbf{x}_i$  are support vectors with non-zero Lagrangian multipliers  $\alpha_i$ .
- Both in the dual formulation of the problem and in the solution, training points appear only inside inner products:

Find  $\alpha_1 \dots \alpha_N$  such that

$\mathbf{Q}(\boldsymbol{\alpha}) = \sum \alpha_i - \frac{1}{2} \sum \sum \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$  is maximized and

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

(2)  $0 \leq \alpha_i \leq C$  for all  $\alpha_i$

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