

# Support Vector Machines

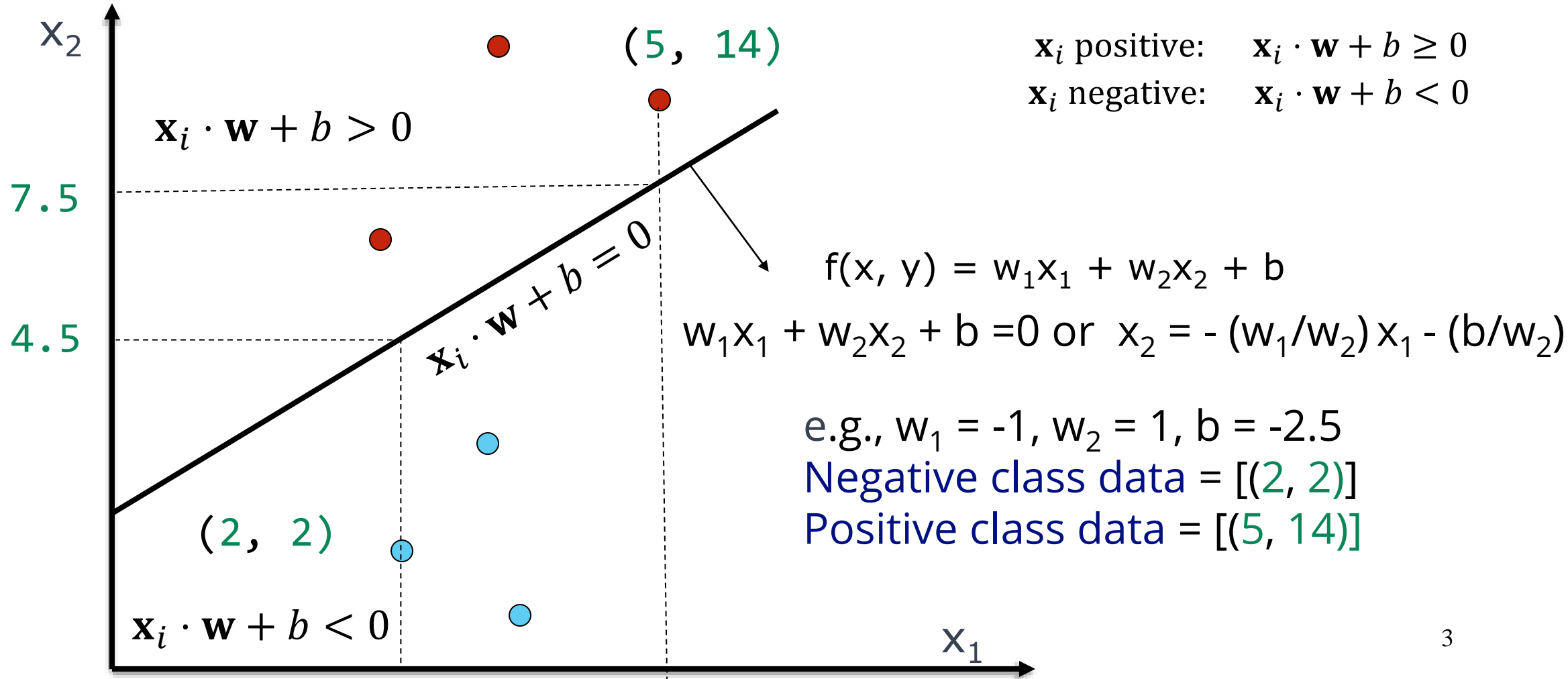
Slides by Svetlana Lazebnik  
(with minor modification)

# Classifiers

- Given a feature representation for images/data, how do we learn a model for distinguishing features from different classes?
- Today:
  - Linear classifiers: support vector machines

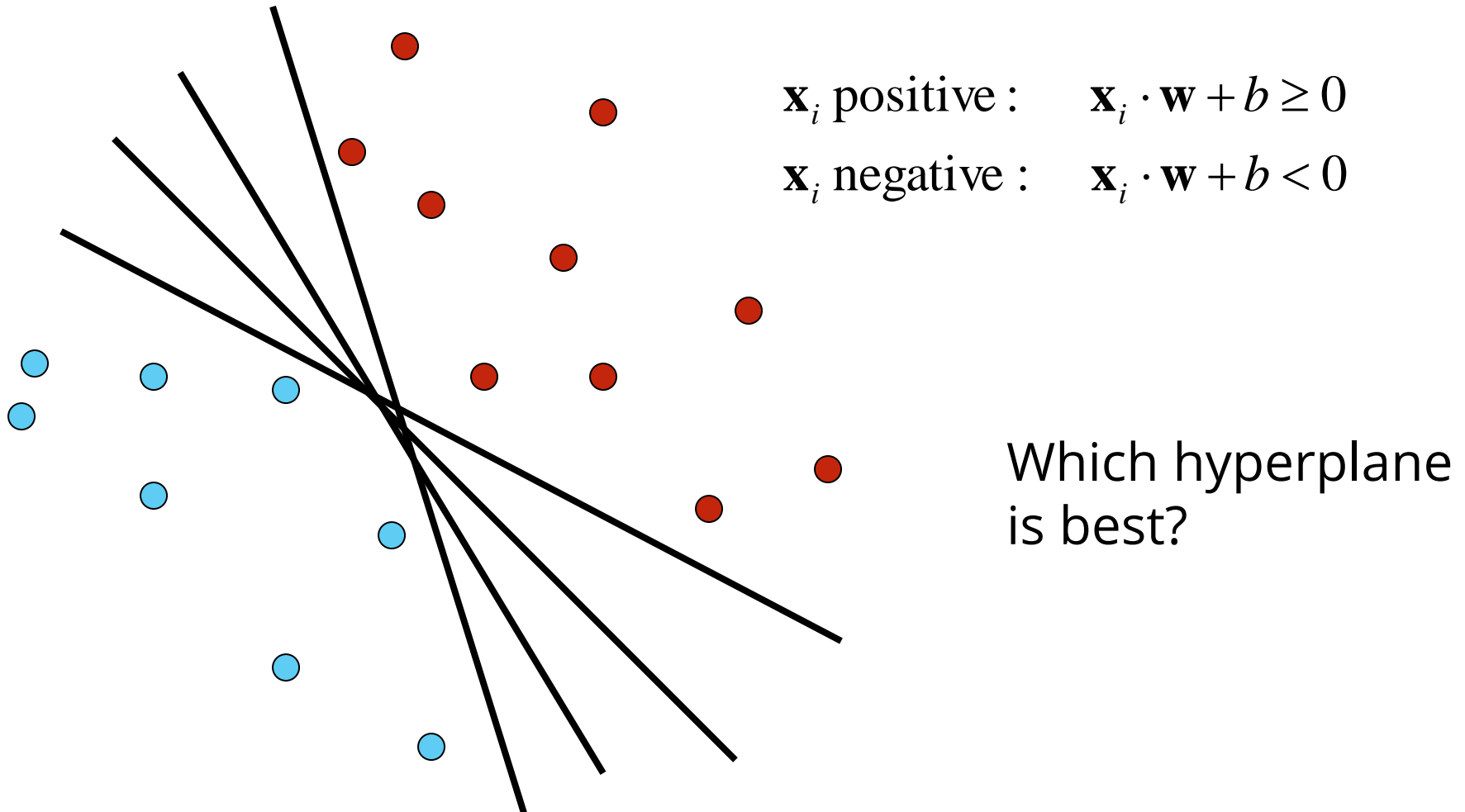
# Linear Classifiers

Find linear function (*hyperplane*) to separate positive and negative examples



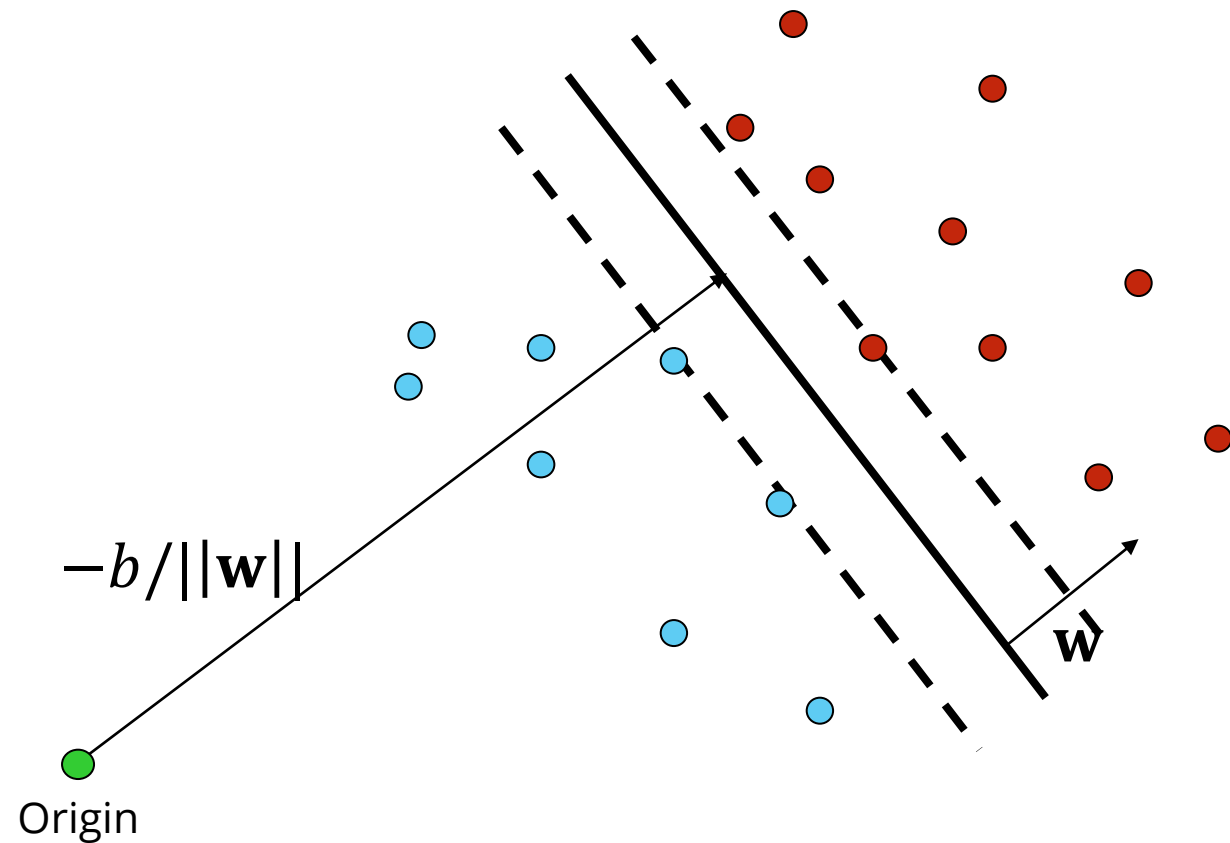
# Linear Classifiers

Find linear function (*hyperplane*) to separate positive and negative examples



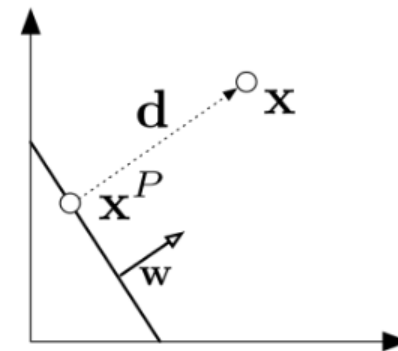
# Support Vector Machines

Find hyperplane that maximizes the *margin* between the positive (say  $y_i=1$ ) and negative (say  $y_i=-1$ ) examples



$$\begin{aligned} \mathbf{x}_i \text{ positive } (y_i = 1)^* : \mathbf{x}_i \cdot \mathbf{w} + b &\geq 1 \\ \mathbf{x}_i \text{ negative } (y_i = -1) : \mathbf{x}_i \cdot \mathbf{w} + b &\leq -1 \\ \Rightarrow y_i (\mathbf{x}_i \cdot \mathbf{w} + b) - 1 &\geq 0 \quad \forall i \end{aligned}$$

distance of a point to the hyperplane



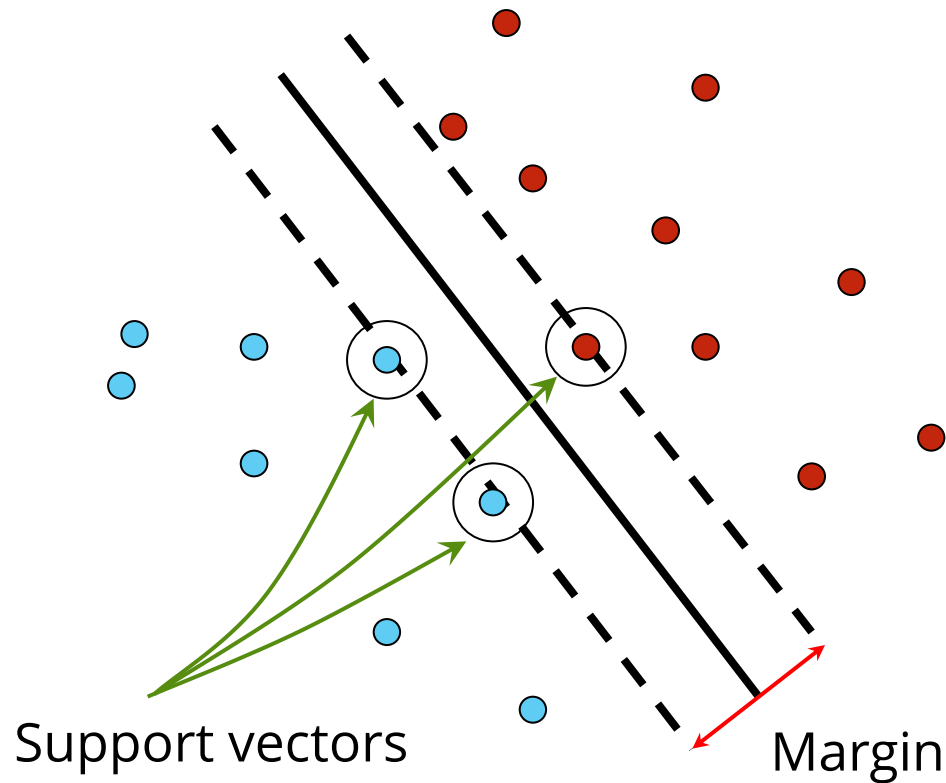
$$\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$$

(for derivation see, e.g., [here](#))

\* we replace 0 with 1 and -1 By setting the thresholds at -1 and 1, it ensures that the data points on the margin, which are called support vectors, have a margin of exactly 1 unit from the decision boundary

# Support Vector Machines

Find hyperplane that maximizes the *margin* between the positive (say  $y_i=1$ ) and negative (say  $y_i=-1$ ) examples



$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

For support vectors,  $\mathbf{x}_i \cdot \mathbf{w} + b = \pm 1$

Distance between point and hyperplane:  $\frac{|\mathbf{x}_i \cdot \mathbf{w} + b|}{\|\mathbf{w}\|}$

Therefore, the margin is  $2 / \|\mathbf{w}\|$  for support vectors

# Finding the Maximum-Margin Hyperplane

1. Maximize margin  $2 / \|\mathbf{w}\|$
2. Correctly classify all training data:

$$\mathbf{x}_i \text{ positive } (y_i = 1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \geq 1$$

$$\mathbf{x}_i \text{ negative } (y_i = -1): \quad \mathbf{x}_i \cdot \mathbf{w} + b \leq -1$$

Quadratic optimization problem:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i (\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$$


Objective-quadratic

constraints are linear

can be solve  
efficiently with  
QCQP  
(Quadratically  
Constrained  
Quadratic  
Program)  
solver

# Finding the Maximum-Margin Hyperplane

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



Learned weight  
(nonzero only for support vectors)

- Support vectors are critical elements of the training set in SVM.
- They are the data points that lie closest to the decision boundary.
- If all other training points were removed or moved within their respective classes, the same separating hyperplane would be found during the training process.



# Finding the Maximum-Margin Hyperplane

- **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\*

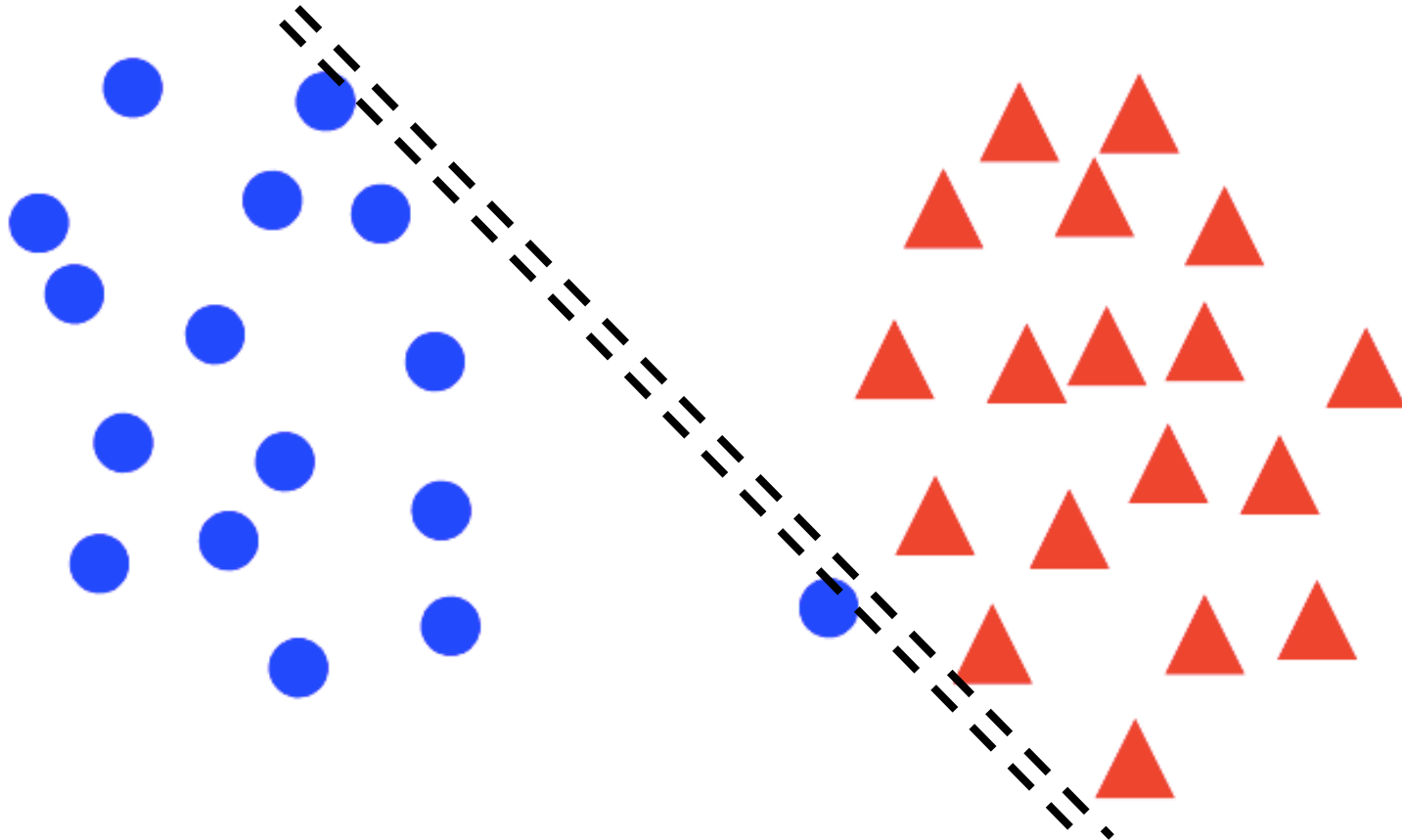
- Classification function (decision boundary):

$$\mathbf{w} \cdot \mathbf{x} + b = \sum_i \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + b$$

- 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

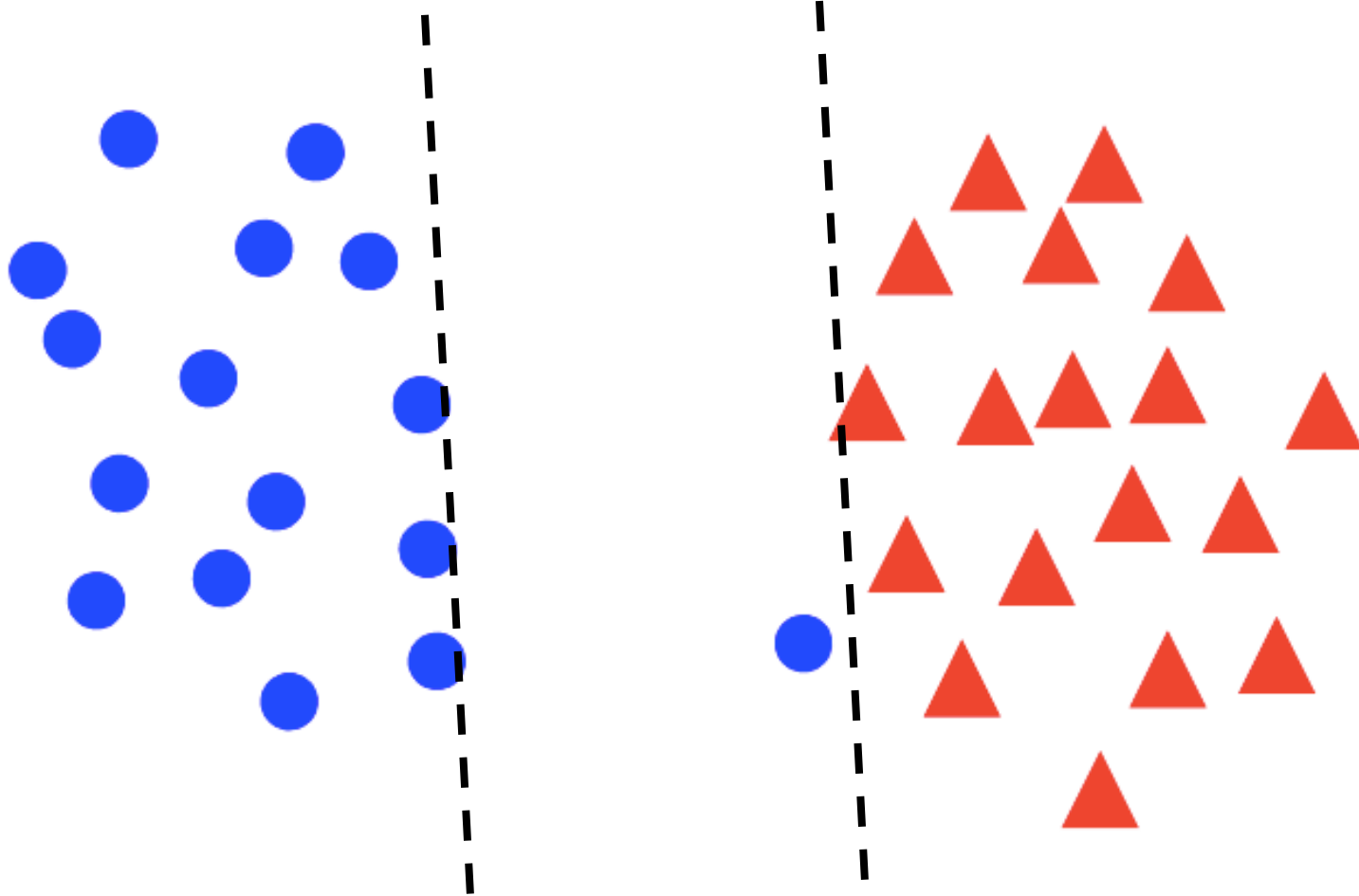
# "Soft margin" formulation

- What about non-separable data?
- And even for separable data, we may prefer a larger margin with a few constraints violated



# "Soft margin" formulation

- What about non-separable data?
- And even for separable data, we may prefer a larger margin with a few constraints violated



# What if the Data is Not Linearly Separable?

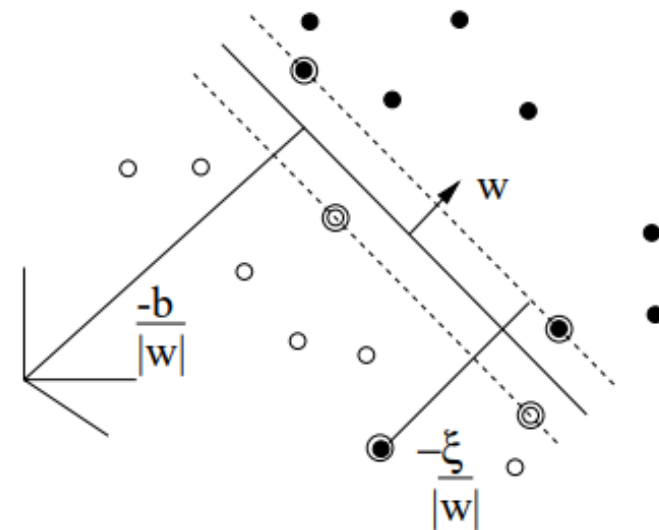
- Separable:  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$  subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

- Idea

Non-separable: change the constrain  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$

Minimize this reduction  $\sum_{i=1}^n \xi_i$

do we need to do this for all?



# What if the Data is Not Linearly Separable?

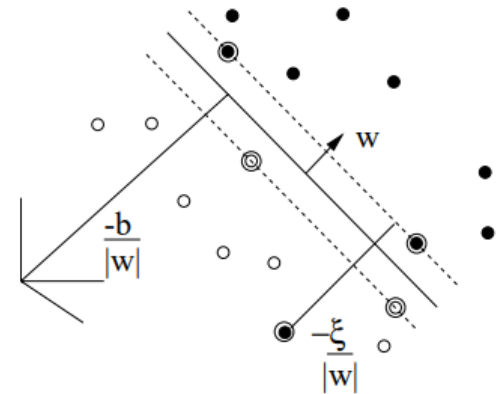
- Separable:  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2$  subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1$

- Non-separable:  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$   
subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i \geq 0$

- $C$ : tradeoff constant,  $\xi_i$ : *slack variable* (positive)

- Whenever margin is  $\geq 1$ ,  $\xi_i = 0$

- Whenever margin is  $< 1$ ,  $\xi_i = 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b)$



# What if the Data is Not Linearly Separable?

Non-separable  $\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i$

subject to  $y_i(\mathbf{w} \cdot \mathbf{x}_i + b) - 1 + \xi_i \geq 0$

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

Maximize  
margin

Minimize classification  
mistakes



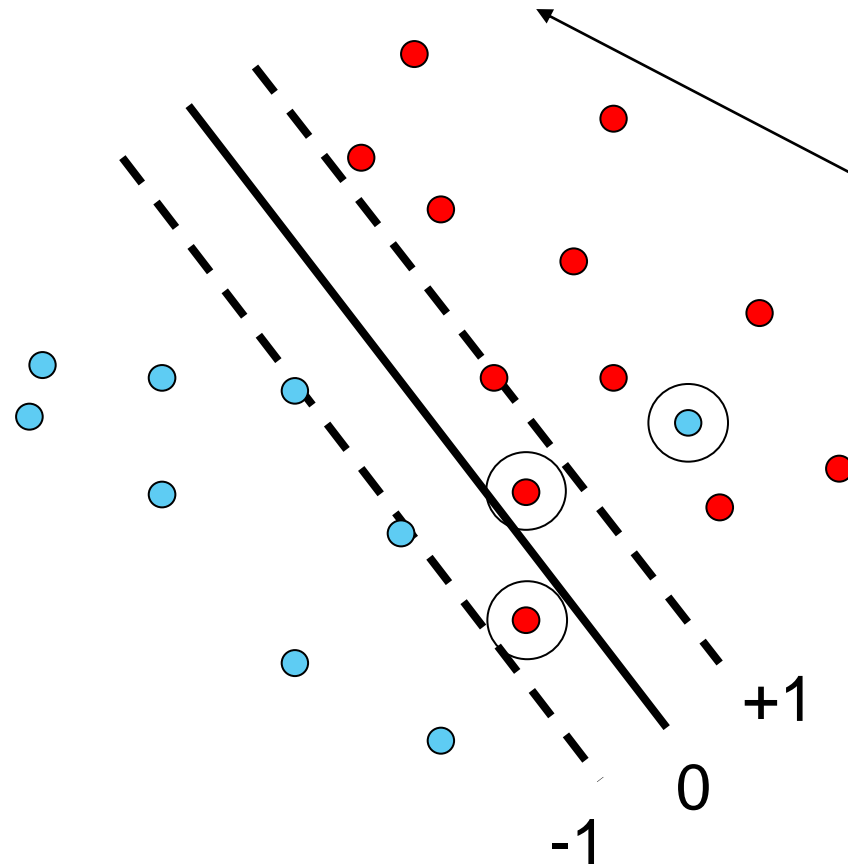
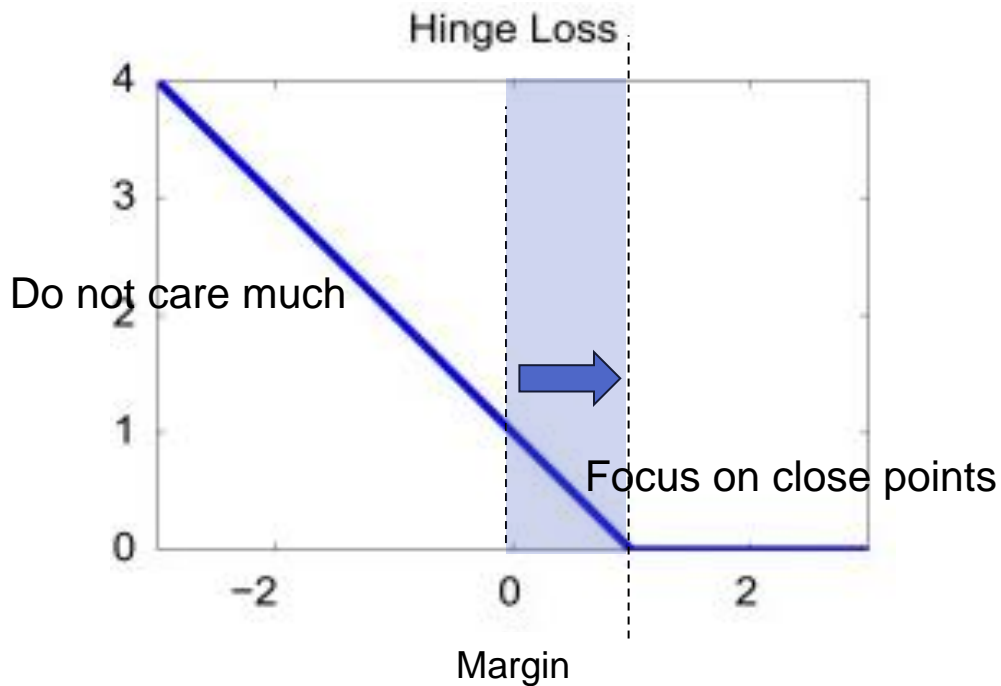
more mistakes

# What if the Data is Not Linearly Separable?

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

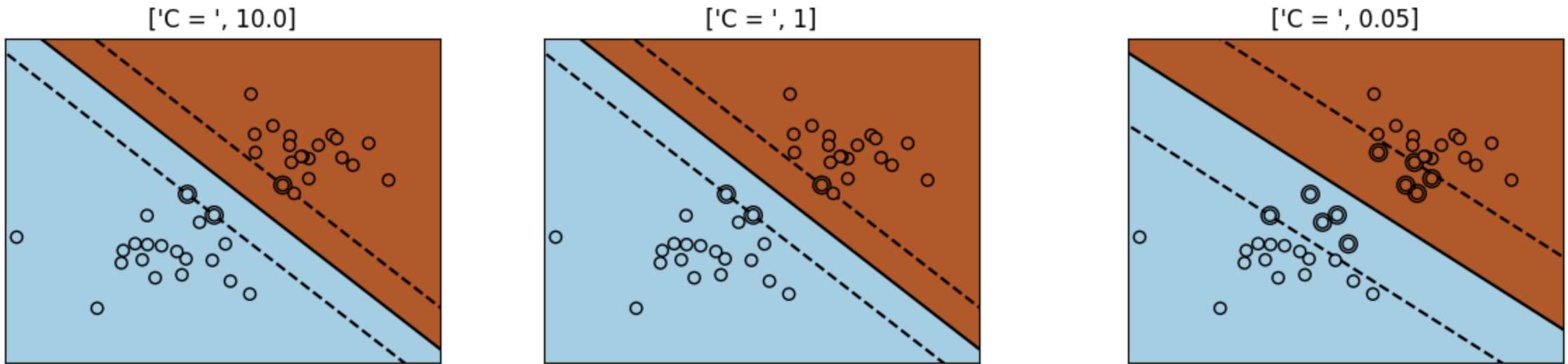
$$\max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$

The hinge loss function is designed to penalize misclassified points and encourage the SVM to maximize the margin between the classes.



# What if the Data is Not Linearly Separable?

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \max(0, 1 - y_i(\mathbf{w} \cdot \mathbf{x}_i + b))$$



“A small value of  $C$  includes more/all the observations, allowing the margins to be calculated using all the data in the area.” [SVM Margins Example — scikit-learn 1.3.0 documentation](#)



# SGD update for SVM

$$l(w, x_i, y_i) = \frac{\lambda}{2n} \|w\|^2 + \max[0, 1 - y_i w^T x_i]$$

$$\nabla l(w, x_i, y_i) = \frac{\lambda}{n} w - \mathbb{I}[y_i w^T x_i < 1] y_i x_i$$

$$\text{Recall: } \frac{d}{da} \max(0, a) = \mathbb{I}[a > 0]$$

# SGD update for SVM

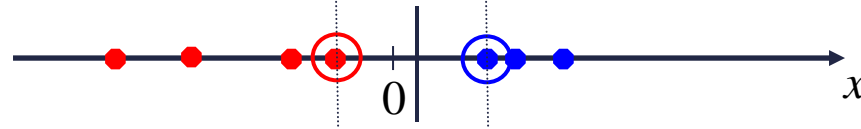
$$l(w, x_i, y_i) = \frac{\lambda}{2n} \|w\|^2 + \max[0, 1 - y_i w^T x_i]$$

$$\nabla l(w, x_i, y_i) = \frac{\lambda}{n} w - \mathbb{I}[y_i w^T x_i < 1] y_i x_i$$

- SGD update:
  - If  $y_i w^T x_i \geq 1$ :  $w \leftarrow w - \eta \frac{\lambda}{n} w$
  - If  $y_i w^T x_i < 1$ :  $w \leftarrow w + \eta \left( y_i x_i - \frac{\lambda}{n} w \right)$

# Nonlinear SVMs

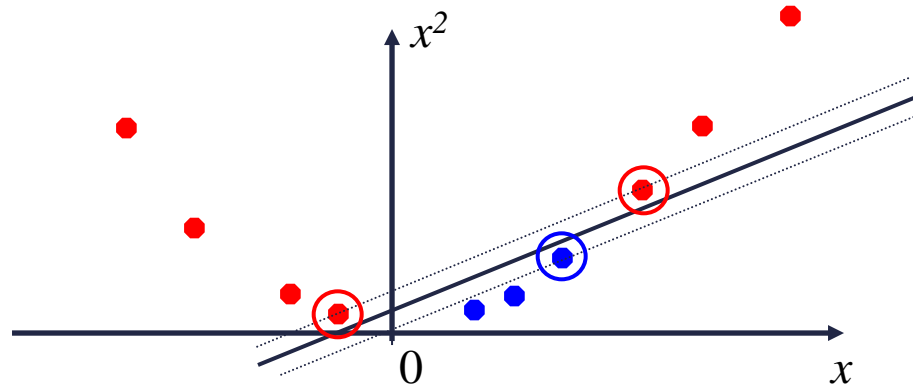
- Datasets that are linearly separable work out great:



- But what if the dataset is just too hard?

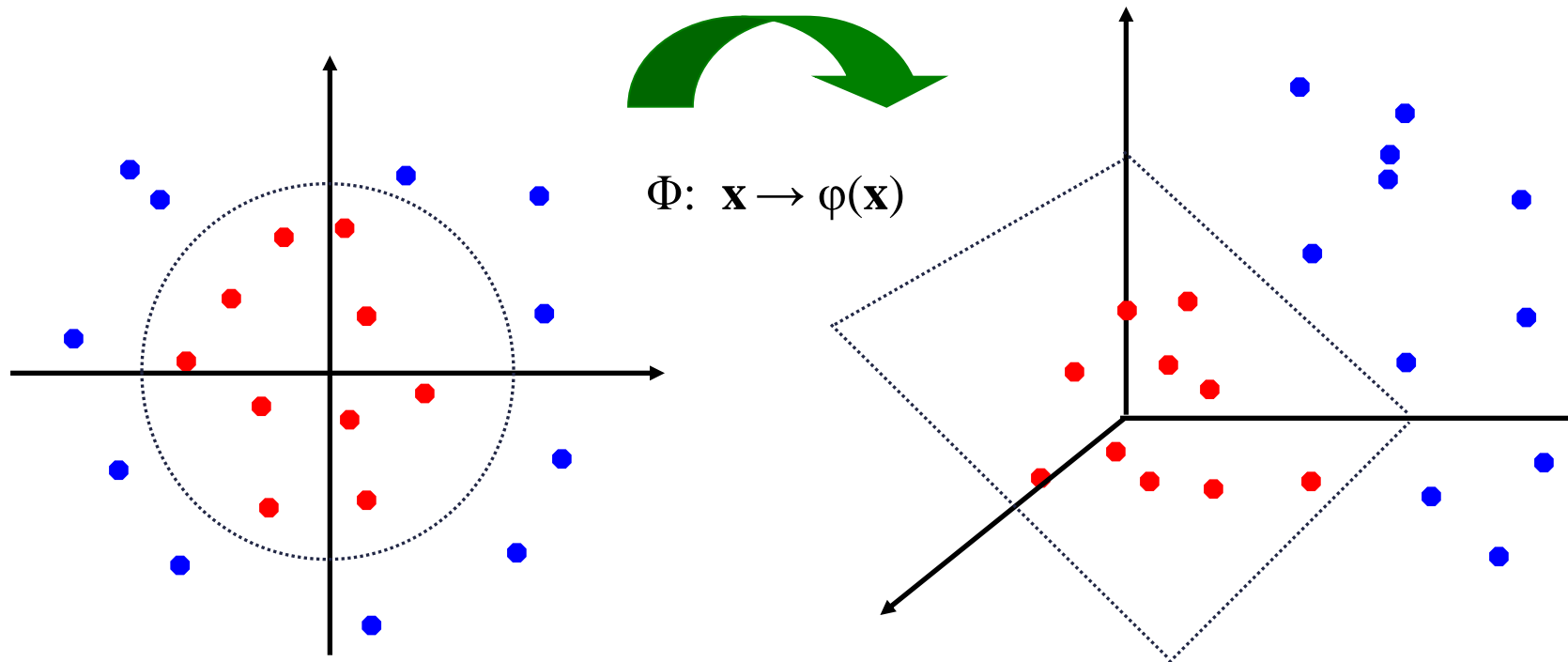


- We can map it to a higher-dimensional space:



# Nonlinear SVMs

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



# Nonlinear SVMs

- *The kernel trick*: instead of explicitly computing the lifting transformation  $\varphi(\mathbf{x})$ , define a kernel function  $K$  such that

$$K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x}) \cdot \varphi(\mathbf{y})$$

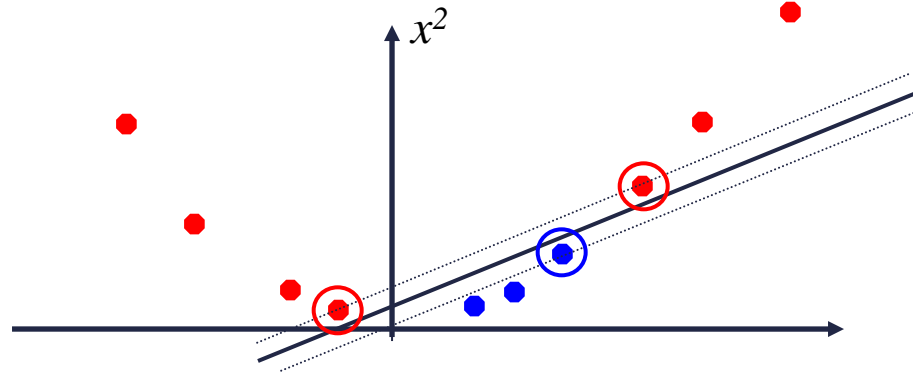
(to be valid, the kernel function must satisfy *Mercer's condition*)

- This gives a nonlinear decision boundary in the original feature space:

$$\sum_i \alpha_i y_i \varphi(\mathbf{x}_i) \cdot \varphi(\mathbf{x}) + b = \sum_i \alpha_i y_i K(\mathbf{x}_i, \mathbf{x}) + b$$

# Nonlinear Kernel: Example

- Consider the mapping  $\varphi(x) = (x, x^2)$

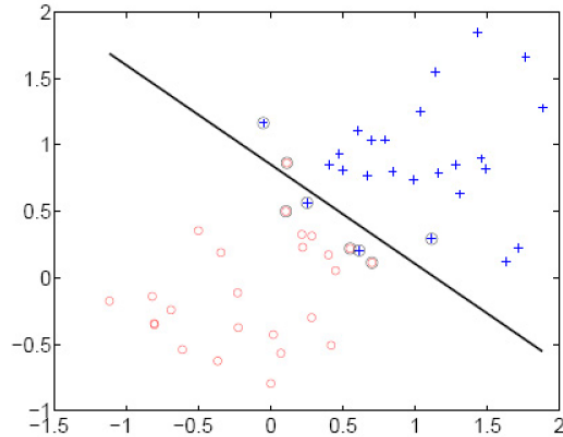


$$\varphi(x) \cdot \varphi(y) = (x, x^2) \cdot (y, y^2) = xy + x^2 y^2$$

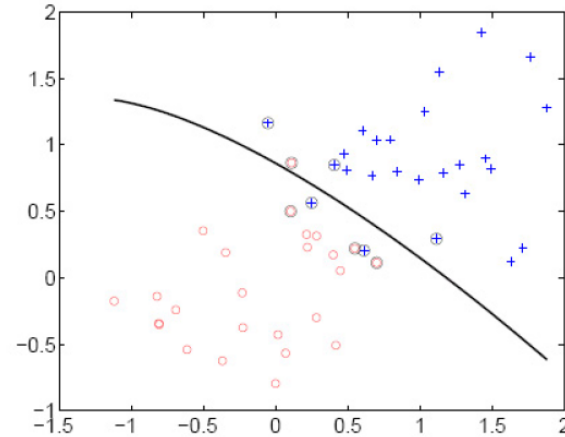
$$K(x, y) = xy + x^2 y^2$$

# Polynomial Kernel:

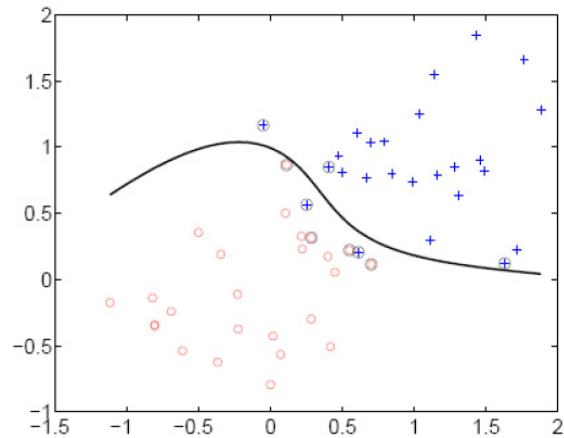
$$K(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x} \cdot \mathbf{y})^d$$



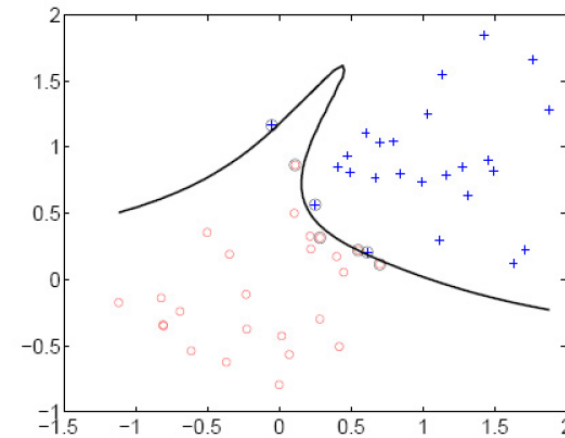
linear



2<sup>nd</sup> order polynomial



4<sup>th</sup> order polynomial



8<sup>th</sup> order polynomial

Higher values of  $d$  allow SVM to fit complex and intricate decision boundaries, but they may also lead to overfitting if not appropriately tuned.

# Gaussian Kernel

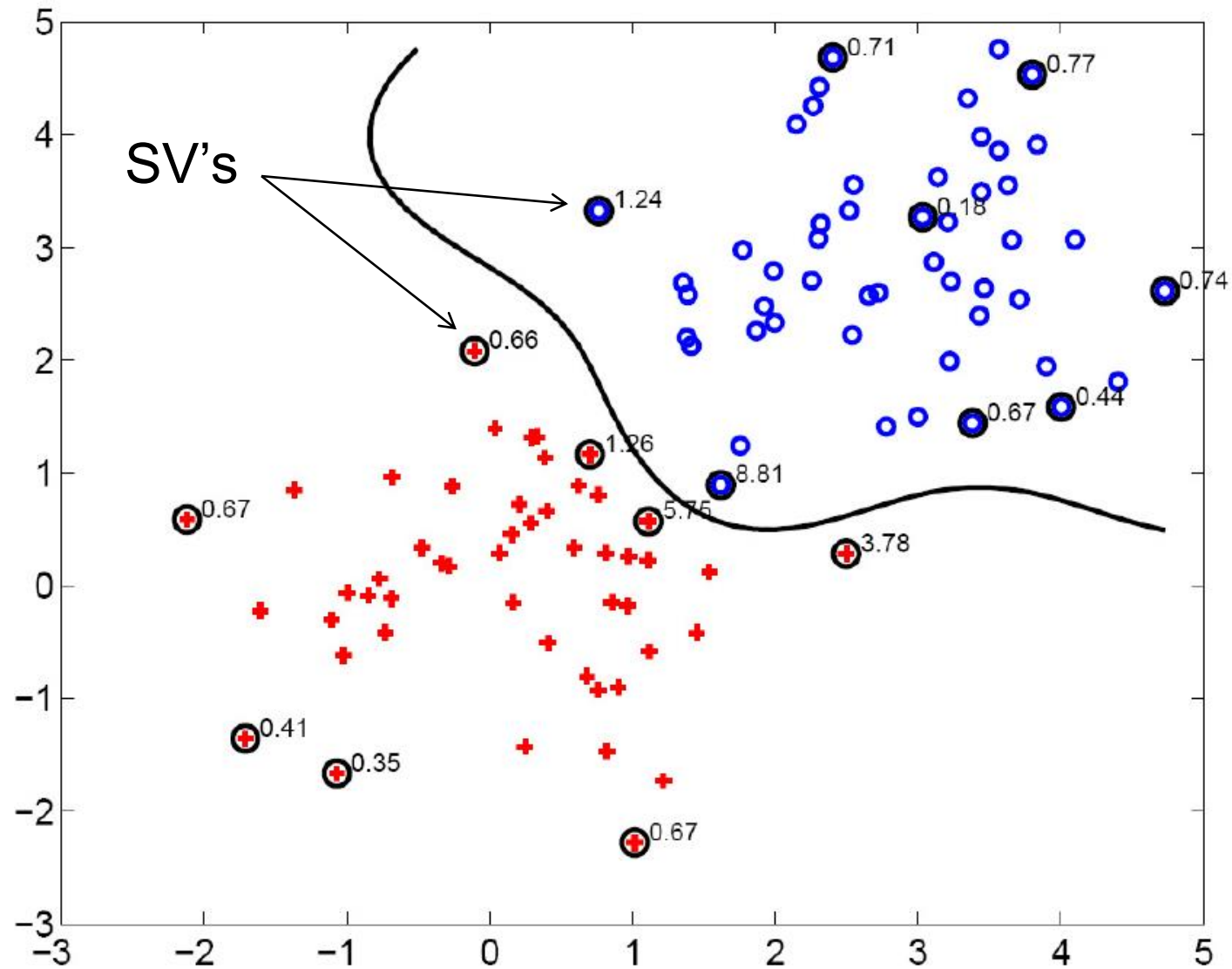
- Also known as the radial basis function (RBF) kernel:

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

- The corresponding mapping  $\phi(\mathbf{x})$  is infinite-dimensional!



# Gaussian Kernel



# Gaussian Kernel

- Also known as the radial basis function (RBF) kernel:

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

- The corresponding mapping  $\phi(\mathbf{x})$  is infinite-dimensional!
- What is the role of parameter  $\sigma$ ?
  - What if  $\sigma$  is close to zero?
  - What if  $\sigma$  is very large?

A small  $\sigma$  can lead to overfitting (The decision boundary will be highly sensitive to individual data points), while a large  $\sigma$  can result in underfitting. Large  $\sigma$  Gaussian kernel becomes more spread out and smoother. The kernel assigns similar weights to a larger region around each data point. A larger  $\sigma$  results in a more generalized decision boundary, and the SVM model is less likely to overfit the training data. However, too large a value of  $\sigma$  may lead to the decision boundary becoming too smooth, potentially underfitting the data and not capturing the intricacies of the data distribution

# Kernels for Bags of Features

- Histogram intersection kernel:

$$I(\mathbf{h}_1, \mathbf{h}_2) = \sum_{i=1}^N \min(\mathbf{h}_1(i), \mathbf{h}_2(i))$$

- Hellinger kernel:  $K(\mathbf{h}_1, \mathbf{h}_2) = \sum_{i=1}^N \sqrt{\mathbf{h}_1(i) \mathbf{h}_2(i)}$

- Generalized Gaussian kernel:

$$K(\mathbf{h}_1, \mathbf{h}_2) = \exp\left(-\frac{1}{A} D(\mathbf{h}_1, \mathbf{h}_2)^2\right)$$

- $D$  can be L1, Euclidean,  $\chi^2$  distance, etc.

# Summary: SVMs for Image Classification

1. Pick an image representation (in our case, bag of features)
2. Pick a kernel function for that representation
3. Compute the matrix of kernel values between every pair of training examples
4. Feed the kernel matrix into your favorite SVM solver to obtain support vectors and weights
5. At test time: compute kernel values for your test example and each support vector, and combine them with the learned weights to get the value of the decision function

# What about Multi-Class SVMs?

- Unfortunately, there is no “definitive” multi-class SVM formulation
- In practice, we have to obtain a multi-class SVM by combining multiple two-class SVMs
- One vs. rest (ovr)
  - Training: learn an SVM for each class vs. the others
  - Testing: apply each SVM to test example and assign to it the class of the SVM that returns the highest decision value
- One vs. one (ovo)
  - Training: learn an SVM for each pair of classes
  - Testing: each learned SVM “votes” for a class to assign to the test example

# SVMs: Pros and Cons

- Pros

- Effective in high dimensional spaces.
- Still effective in cases where number of dimensions is greater than the number of samples.
- Uses a subset of training points in the decision function (called support vectors), so it is also memory efficient.
- Many publicly available SVM packages:  
<http://www.kernel-machines.org/software>
- Kernel-based framework is very powerful, flexible.
- SVMs work very well in practice, even with very small training sample sizes

- Cons

- No “direct” multi-class SVM, must combine two-class SVMs.
- If the number of features is much greater than the number of samples, avoid over-fitting in choosing Kernel functions and regularization term is crucial.
- SVMs do not directly provide probability estimates.
- Computation, memory
  - During training time, must compute matrix of kernel values for every pair of examples
  - Learning can take a very long time for large-scale problems

# SVMs for Large-Scale Datasets

- Efficient *linear* solvers
  - [LIBLINEAR](#), [PEGASOS](#)
- Explicit approximate embeddings: define an explicit mapping  $\phi(\mathbf{x})$  such that  $\phi(\mathbf{x}) \cdot \phi(\mathbf{y})$  approximates  $K(\mathbf{x}, \mathbf{y})$  and train a linear SVM on top of that embedding
  - Random Fourier features for the Gaussian kernel (Rahimi and Recht, 2007)
  - Embeddings for additive kernels, e.g., histogram intersection (Maji et al., 2013, Vedaldi and Zisserman, 2012)

# Summary: Classifiers

- Nearest-neighbor and k-nearest-neighbor classifiers
- Support vector machines
  - Linear classifiers
  - Margin maximization
  - Non-separable case
  - The kernel trick
  - Multi-class SVMs
  - Large-scale SVMs
- Of course, there are many other classifiers out there
  - Neural networks, boosting, decision trees/forests, ...