

# ITCS 5356: Intro to Machine Learning

---

## The Average Perceptron The Kernel Trick

Razvan C. Bunescu

Department of Computer Science @ CCI

[razvan.bunescu@charlotte.edu](mailto:razvan.bunescu@charlotte.edu)

# Linear Discriminant Functions: Two classes ( $K = 2$ )

---

- Use a linear function of the input vector:

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

*weight vector*

*bias = - threshold*

- Decision:

$\mathbf{x} \in C_1$  if  $h(\mathbf{x}) \geq 0$ , otherwise  $\mathbf{x} \in C_2$ .

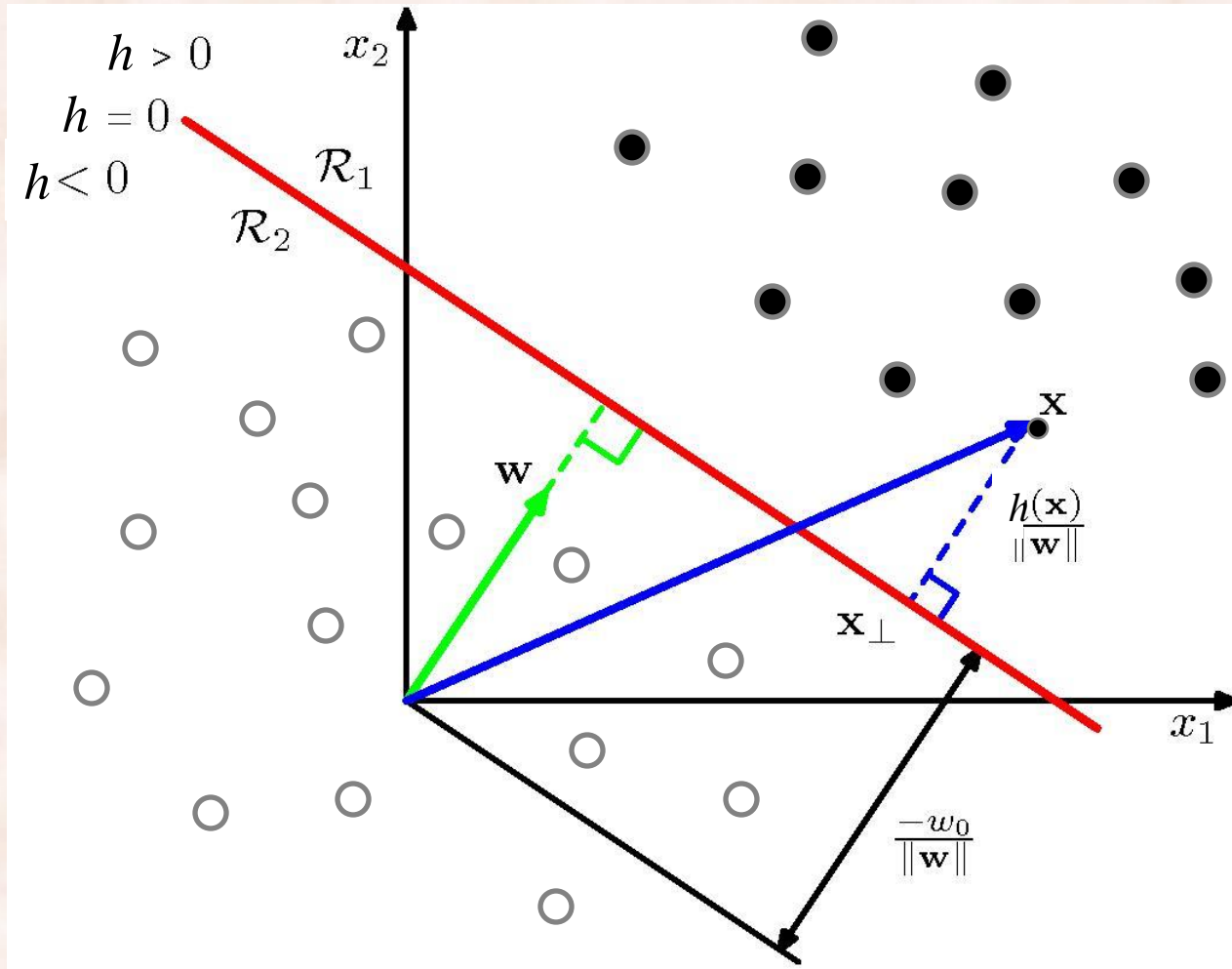
$\Rightarrow$  **decision boundary** is hyperplane  $h(\mathbf{x}) = 0$ .

- Properties:

- $\mathbf{w}$  is orthogonal to vectors lying within the decision surface.
- $w_0$  controls the location of the decision hyperplane.

# Geometric Interpretation

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0 = w_1 x_1 + w_2 x_2 + w_0$$



# Linear Models for Classification

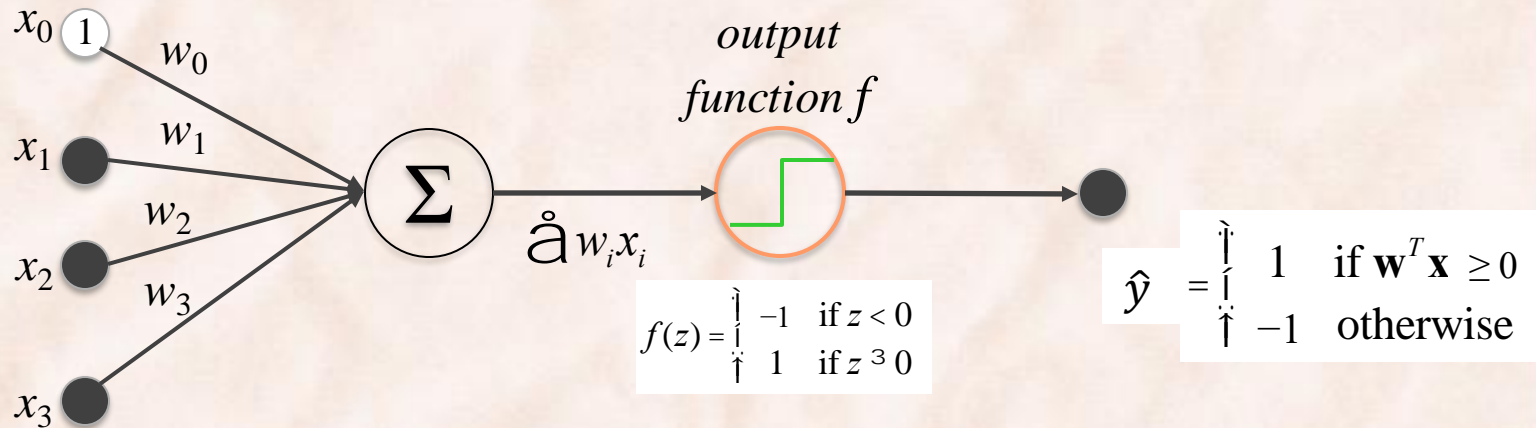
---

- We want to use a linear function of the feature vector:

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- How to find  $\mathbf{w}$  automatically? Use ML!
  - **Perceptron.**
  - **Logistic Regression.**
- What if the data is not linearly separable? Make it!
  - Engineer new features (LR) or use kernels (Perceptron).
  - Learn new features (**Neural Networks**).

# Linear Discriminant Classification



- Assume classes  $T = \{c_1, c_2\} = \{1, -1\}$ .
- Training set is  $(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_n, y_n)$ .

$$\mathbf{x} = [1, x_1, x_2, \dots, x_k]^T$$

$$\hat{y} = \text{sgn}(\mathbf{w}^T \mathbf{x}) = \text{sgn}(w_0 + w_1 x_1 + \dots + w_k x_k)$$

*a linear discriminant function*

# Linear Discriminant Classification: Objective Function

---

- Learning = finding the “right” parameters  $\mathbf{w}^T = [w_0, w_1, \dots, w_k]$ 
  - Find  $\mathbf{w}$  that minimizes an *error function*  $E(\mathbf{w})$  which measures the misfit between  $\hat{t}(\mathbf{x}_n)$  and  $t_n$ .
- **Least Squares** error function?

$$J(\mathbf{w}) = \frac{1}{2N} \sum_{n=1}^N (\hat{y}_n - y_n)^2$$

$$\hat{y} = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq 0 \\ -1 & \text{otherwise} \end{cases}$$

4 times # of mistakes

# Least Squares vs. Perceptron Criterion

---

- **Least Squares**  $\Rightarrow$  cost is # of misclassified patterns:
  - Piecewise constant function of  $\mathbf{w}$  with discontinuities.
  - Cannot find closed form solution for  $\mathbf{w}$  that minimizes cost.
  - Cannot use gradient methods (gradient zero almost everywhere).
- **Perceptron Criterion:**
  - Set labels to be +1 and -1. Want  $\mathbf{w}^T \mathbf{x}_n > 0$  for  $y_n = 1$ , and  $\mathbf{w}^T \mathbf{x}_n < 0$  for  $y_n = -1$ .
    - $\Rightarrow$  would like to have  $\mathbf{w}^T \mathbf{x}_n y_n > 0$  for all patterns.
    - $\Rightarrow$  want to minimize  $-\mathbf{w}^T \mathbf{x}_n y_n$  for all misclassified patterns  $M$ .

$$\Rightarrow \text{minimize } E_p(\mathbf{w}) = -\sum_{n \in M} \mathbf{w}^T \mathbf{x}_n y_n$$

# Stochastic Gradient Descent

---

- **Perceptron Criterion:**

$$\text{minimize } E_p(\mathbf{w}) = -\sum_{n \in M} \mathbf{w}^T \mathbf{x}_n y_n$$

- Update parameters  $\mathbf{w}$  sequentially **after each mistake:**

$$\begin{aligned} \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - h \tilde{\nabla} E_p(\mathbf{w}^{(t)}, \mathbf{x}_n) \\ &= \mathbf{w}^{(\tau)} + \eta \mathbf{x}_n y_n \end{aligned}$$

- The magnitude of  $\mathbf{w}$  is inconsequential  $\Rightarrow$  can set  $\eta = 1$ .

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} + \mathbf{x}_n y_n$$

Prove it.



# The Perceptron Algorithm: Two Classes

---

1. **initialize** parameters  $\mathbf{w} = 0$

2. **for**  $n = 1 \dots N$

3.  $h_n = \mathbf{w}^T \mathbf{x}_n$

4. **if**  $h_n \geq 0$  and  $y_n = -1$

5.  $\mathbf{w} = \mathbf{w} - \mathbf{x}_n$

6. **if**  $h_n \leq 0$  and  $y_n = +1$

7.  $\mathbf{w} = \mathbf{w} + \mathbf{x}_n$

Repeat:

a) until convergence.

b) for a number of epochs  $E$ .

What is the impact of the perceptron update on the score  $\mathbf{w}^T \mathbf{x}_n$  of the misclassified example  $\mathbf{x}_n$ ?

# The Perceptron Algorithm: Two Classes

---

1. **initialize** parameters  $\mathbf{w} = 0$
2. **for**  $n = 1 \dots N$
3.  $h_n = \mathbf{w}^T \mathbf{x}_n$
4. **if**  $h_n y_n \leq 0$  **then**
5.  $\mathbf{w} = \mathbf{w} + y_n \mathbf{x}_n$

Repeat:

- a) until convergence.
- b) for a number of epochs  $E$ .

Loop invariant:  $\mathbf{w}$  is a weighted sum of training vectors:

$$\mathbf{w} = \sum_n \alpha_n y_n \mathbf{x}_n \quad \Rightarrow \quad \mathbf{w}^T \mathbf{x} = \sum_n \alpha_n y_n \mathbf{x}_n^T \mathbf{x}$$

# The Perceptron Algorithm: Two Classes

1. **initialize** parameters  $\mathbf{w} = 0$
2. **for**  $n = 1 \dots N$
3.  $\hat{y}_n = \text{sgn}(\mathbf{w}^T \mathbf{x}_n)$
4. **if**  $\hat{y}_n \neq y_n$  **then**
5.  $\mathbf{w} = \mathbf{w} + y_n \mathbf{x}_n$

$$\text{sgn}(h) = \begin{cases} +1 & \text{if } h > 0, \\ 0 & \text{if } h = 0, \\ -1 & \text{if } h < 0 \end{cases}$$

Repeat:

- a) until convergence.
- b) for a number of epochs  $E$ .

Theorem [[Rosenblatt, 1962](#)]:

If the training dataset is linearly separable, the perceptron learning algorithm is guaranteed to find a solution in a finite number of steps.

- see Theorem 1 (Block, Novikoff) in [[Freund & Schapire, 1999](#)].

# The Perceptron Algorithm: Two Classes

---

**initialize** parameters  $\mathbf{w} = 0$

**for** epoch  $e = 1 \dots E$

mistakes = 0

**for** example  $n = 1 \dots N$

$$\hat{y}_n = \text{sgn}(\mathbf{w}^T \mathbf{x}_n)$$

**if**  $\hat{y}_n \neq y_n$  **then**

$$\mathbf{w} = \mathbf{w} + y_n \mathbf{x}_n$$

mistakes = mistakes + 1

**if** mistakes = 0

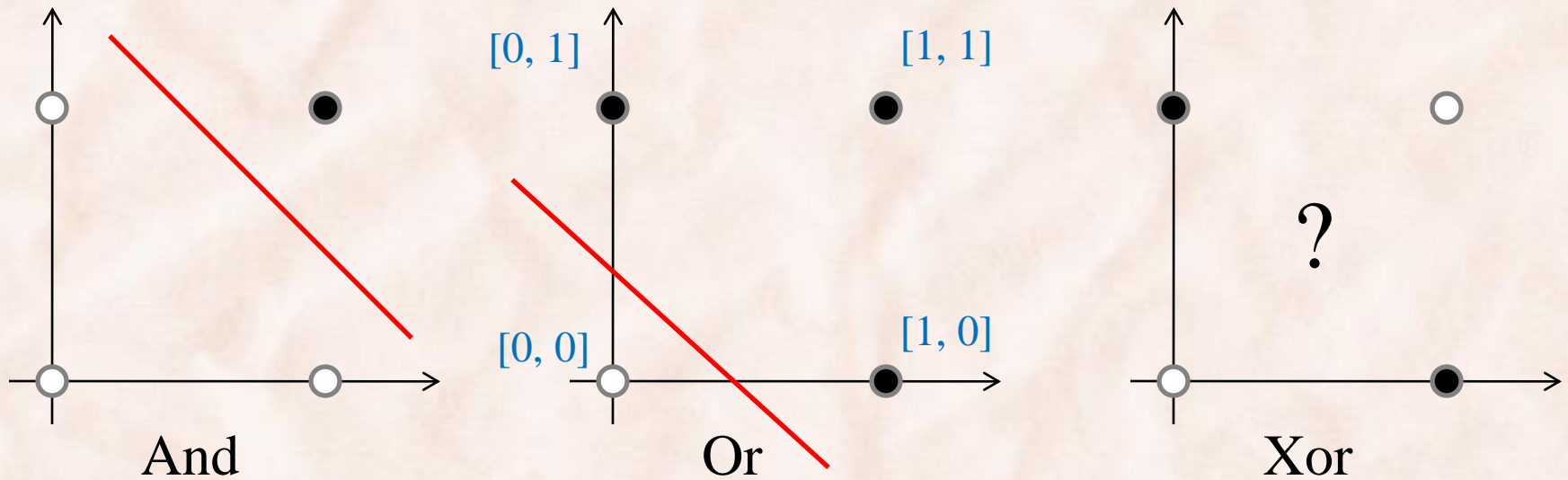
**break** Converged!

$$\text{sgn}(h) = \begin{cases} +1 & \text{if } h > 0, \\ 0 & \text{if } h = 0, \\ -1 & \text{if } h < 0 \end{cases}$$

1 epoch = one pass over all training examples.

# Linear vs. Non-linear Decision Boundaries

---



$$\left. \begin{array}{l} \mathbf{x} = [1, x_1, x_2]^T \\ \mathbf{w} = [w_0, w_1, w_2]^T \end{array} \right\} \Rightarrow \mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + w_0$$

# How to Find Non-linear Decision Boundaries

---

1) Perceptron with manually engineered features:

- Quadratic features.

2) **Kernel methods with non-linear kernels:**

- Quadratic kernels, Gaussian kernels.

*Deep Learning*

3) Self-supervised feature learning (e.g. auto-encoders):

- Plug learned features in any linear classifier.

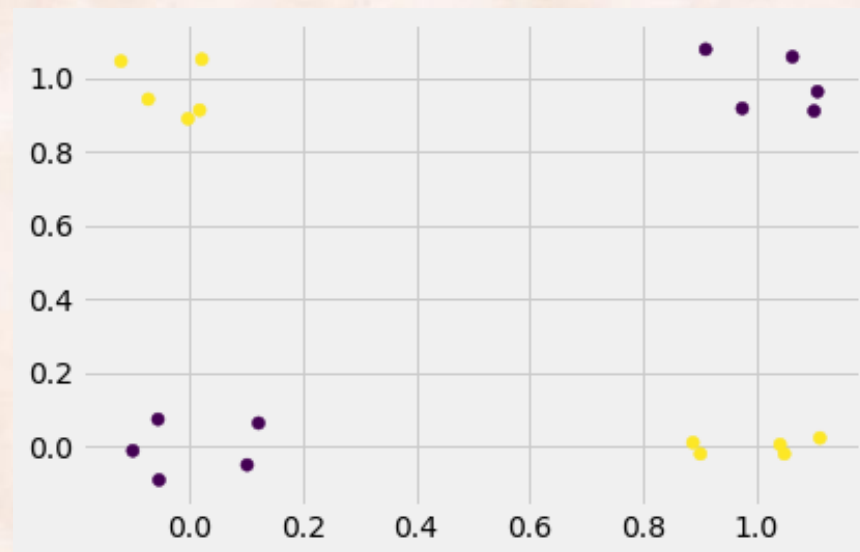
4) Neural Networks with one or more hidden layers:

- Automatically learned features.

# Non-Linear Classification: XOR Dataset

---

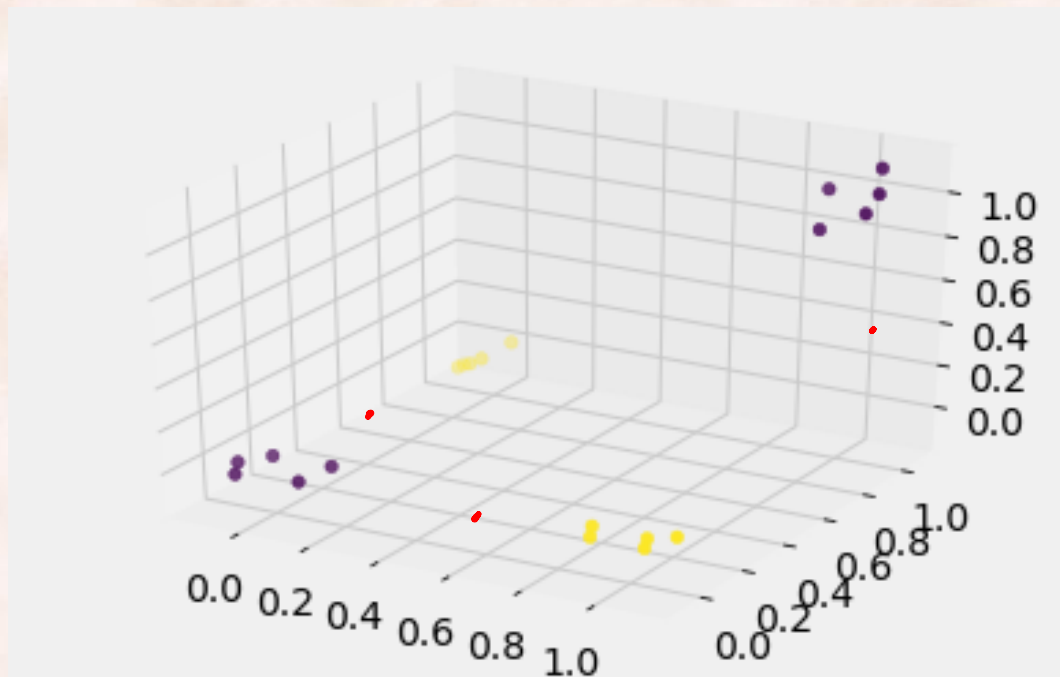
$$\mathbf{x} = [x_1, x_2]$$



# 1) Manually Engineered Features: Add $x_1x_2$

---

$$\mathbf{x} = [x_1, x_2, x_1x_2]$$

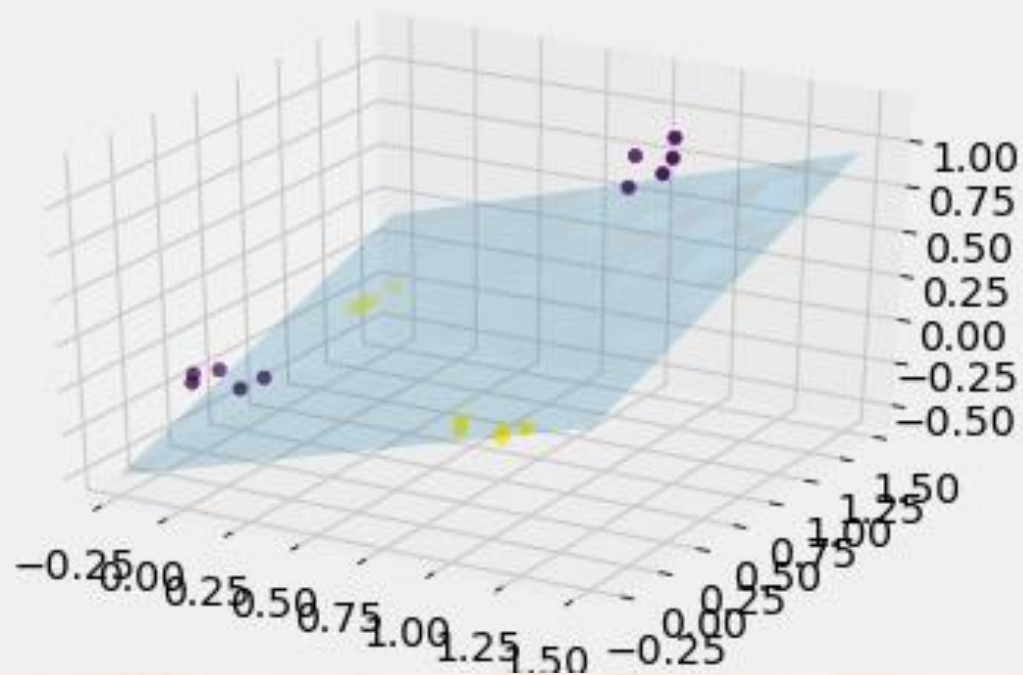




# Logistic Regression with Manually Engineered Features

---

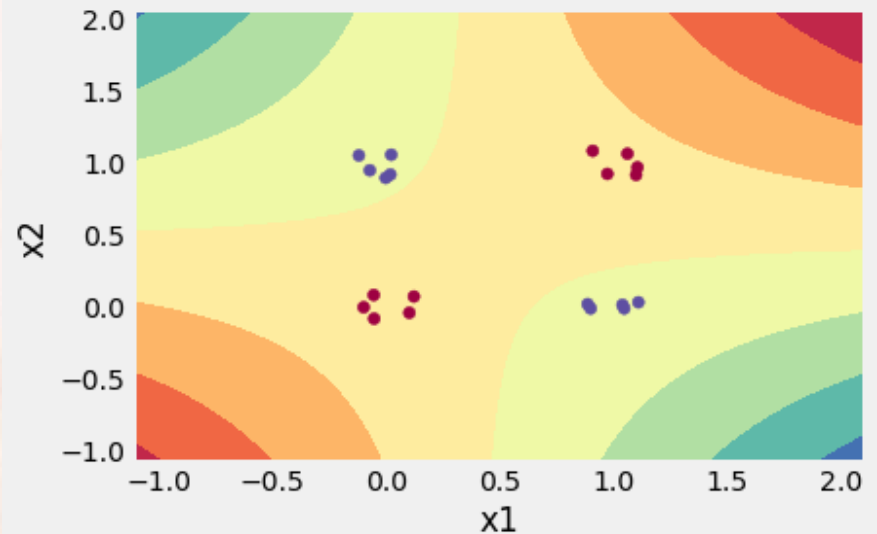
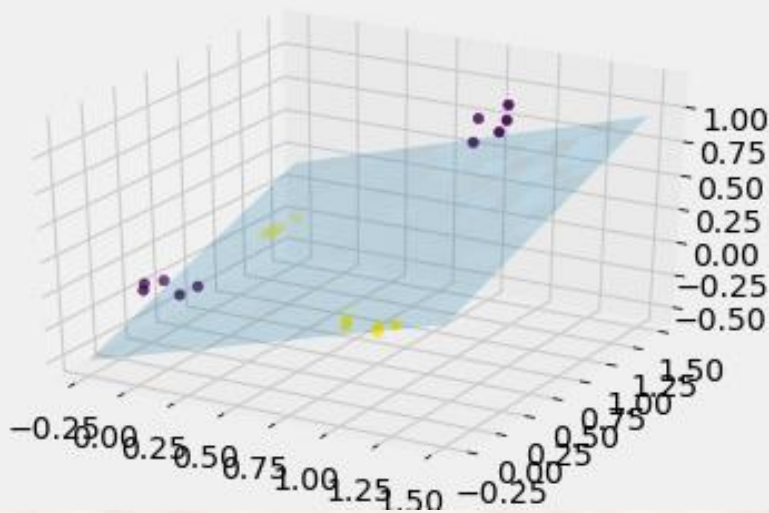
$$\mathbf{x} = [x_1, x_2, x_1x_2]$$



# Perceptron with Manually Engineered Features

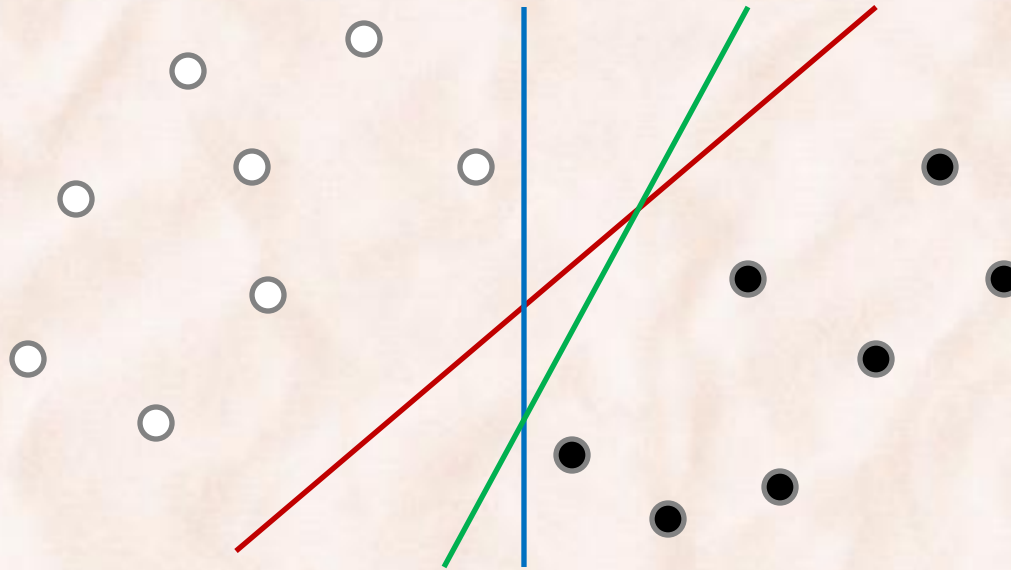
---

Project  $\mathbf{x} = [x_1, x_2, x_1x_2]$  and decision hyperplane back to  $\mathbf{x} = [x_1, x_2]$



# Classifiers & Margin

---



- Which classifier has the smallest generalization error?
  - The one that maximizes the margin [[Computational Learning Theory](#)]
    - **margin** = the distance between the decision boundary and the closest sample.

# Averaged Perceptron: Two Classes

1. **initialize** parameters  $\mathbf{w} = \mathbf{0}$ ,  $\tau = 1$ ,  $\bar{\mathbf{w}} = \mathbf{0}$
2. **for**  $n = 1 \dots N$
3.      $\hat{y}_n = \text{sgn}(\mathbf{w}^T \mathbf{x}_n)$
4.     **if**  $\hat{y}_n \neq y_n$  **then**
5.          $\mathbf{w} = \mathbf{w} + y_n \mathbf{x}_n$
6.          $\bar{\mathbf{w}} = \bar{\mathbf{w}} + \mathbf{w}$
7.          $\tau = \tau + 1$
8. **return**  $\bar{\mathbf{w}}/\tau$

$$\text{sgn}(z) = \begin{cases} +1 & \text{if } z > 0, \\ 0 & \text{if } z = 0, \\ -1 & \text{if } z < 0 \end{cases}$$

Repeat:

- a) until convergence.
- b) for a number of epochs E.

During testing:  $\hat{y} = \text{sgn}(\bar{\mathbf{w}}^T \mathbf{x})$

## 2) Kernel Methods with Non-Linear Kernels

---

- Perceptrons, SVMs can be ‘*kernelized*’:
  1. Re-write the algorithm such that during training and testing feature vectors  $\mathbf{x}$ ,  $\mathbf{y}$  appear only in dot-products  $\mathbf{x}^T \mathbf{y}$ .
  2. Replace dot-products  $\mathbf{x}^T \mathbf{y}$  with *non-linear kernels*  $K(\mathbf{x}, \mathbf{y})$ :
    - $K$  is a kernel if and only if  $\exists \varphi$  such that  $K(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^T \varphi(\mathbf{y})$ 
      - $\varphi$  can be in a much higher dimensional space.
        - » e.g. combinations of up to  $k$  original features
      - $\varphi(\mathbf{x})^T \varphi(\mathbf{y})$  can be computed efficiently without enumerating  $\varphi(\mathbf{x})$  or  $\varphi(\mathbf{y})$ .

# The Perceptron Representer Theorem

---

1. **initialize** parameters  $\mathbf{w} = 0$
2. **for**  $n = 1 \dots N$
3.      $\hat{y}_n = \text{sgn}(\mathbf{w}^T \mathbf{x}_n)$
4.     **if**  $\hat{y}_n \neq y_n$  **then**
5.          $\mathbf{w} = \mathbf{w} + y_n \mathbf{x}_n$

Repeat:

- a) until convergence.
- b) for a number of epochs  $E$ .

Loop invariant:  $\mathbf{w}$  is a weighted sum of training vectors:

$$\mathbf{w} = \sum_{n=1..N} \alpha_n y_n \mathbf{x}_n \Rightarrow \mathbf{w}^T \mathbf{x} = \sum_{n=1..N} \alpha_n y_n \mathbf{x}_n^T \mathbf{x}$$

# Kernel Perceptron: Two Classes

---

define  $\mathbf{w} = \sum_{n=1..N} \alpha_n y_n \mathbf{x}_n$

1. **initialize** dual parameters  $\alpha_n = 0$

2. **for**  $n = 1 \dots N$

3.  $\hat{y}_n = \text{sgn}(\mathbf{w}^T \mathbf{x}_n)$

4. **if**  $\hat{y}_n \neq y_n$  **then**

5.  $\alpha_n = \alpha_n + 1$

Repeat:

a) until convergence.

b) for a number of epochs E.

Inference:  $\hat{y} = \text{sgn}(\mathbf{w}^T \mathbf{x})$

where  $\mathbf{w}^T \mathbf{x} = \sum_{j=1..N} \alpha_j y_j \mathbf{x}_j^T \mathbf{x} = \sum_{j=1..N} \alpha_j y_j K(\mathbf{x}_j, \mathbf{x})$

# Kernel Perceptron: Two Classes

define  $h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_{j=1..N} \alpha_j y_j \mathbf{x}_j^T \mathbf{x} = \sum_{j=1..N} \alpha_j y_j K(\mathbf{x}_j, \mathbf{x})$

1. **initialize** dual parameters  $\alpha_n = 0$

$$K(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$$

2. **for**  $n = 1 \dots N$

3.  $\hat{y}_n = \text{sgn } f(\mathbf{x}_n)$

4. **if**  $\hat{y}_n \neq y_n$  **then**

5.  $\alpha_n = \alpha_n + 1$

Repeat:

a) until convergence.

b) for a number of epochs E.

Let  $S = \{j | \alpha_j \neq 0\}$  be the set of *support vectors*. Then  $h(\mathbf{x}) = \sum_{j \in S} \alpha_j y_j K(\mathbf{x}_j, \mathbf{x})$

During testing:  $\hat{y} = \text{sgn } h(\mathbf{x})$

$$\mathbf{w} = \sum_{j \in S} \alpha_j y_j \phi(\mathbf{x}_j)$$



# Kernel Perceptron: Complete Pseudocode

**initialize** dual parameters  $\alpha = 0$

**for** epoch  $e = 1 \dots E$

mistakes = 0

**for** example  $n = 1 \dots N$

$$\hat{y}_n = \text{sgn } h(\mathbf{x}_n)$$

**if**  $\hat{y}_n \neq y_n$  **then**

$$\alpha_n = \alpha_n + 1$$

$$\text{mistakes} = \text{mistakes} + 1$$

**if** mistakes = 0

break

**return**  $\alpha$

$$\text{sgn}(h) = \begin{cases} +1 & \text{if } h > 0, \\ 0 & \text{if } h = 0, \\ -1 & \text{if } h < 0 \end{cases}$$

1 epoch = one pass over all training examples.

$$h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_j \alpha_j y_j \mathbf{x}_j^T \mathbf{x} = \sum_j \alpha_j y_j K(\mathbf{x}_j, \mathbf{x})$$

# Kernel Perceptron: Alternative Formulation

$$\text{define } h(\mathbf{x}) = \mathbf{w}^T \mathbf{x} = \sum_j \alpha_j \mathbf{x}_j^T \mathbf{x} = \sum_j \alpha_j K(\mathbf{x}_j, \mathbf{x})$$

1. **initialize** dual parameters  $\alpha_n = 0$

no  $y_j$  anymore

2. **for**  $n = 1 \dots N$

3.  $\hat{y}_n = \text{sgn } h(\mathbf{x}_n)$

4. **if**  $\hat{y}_n \neq y_n$  **then**

5.  $\alpha_n = \alpha_n + y_n$

Repeat:

a) until convergence.

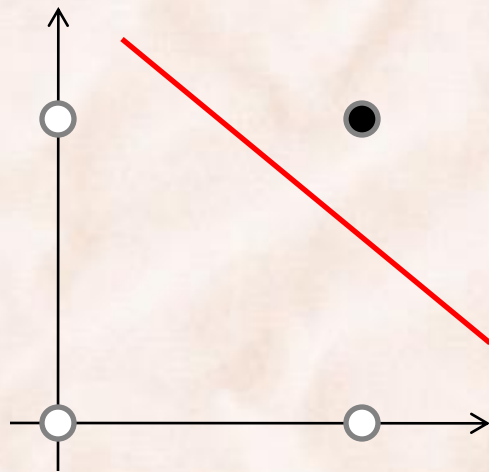
b) for a number of epochs E.

add  $y_n$  instead of 1

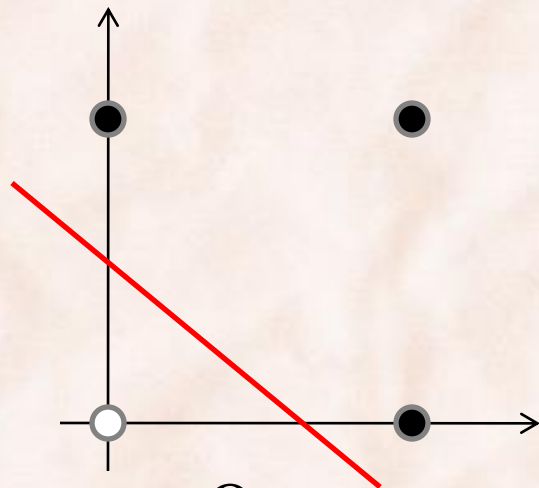
During testing:  $\hat{y} = \text{sgn } h(\mathbf{x})$

# The Perceptron vs. Boolean Functions

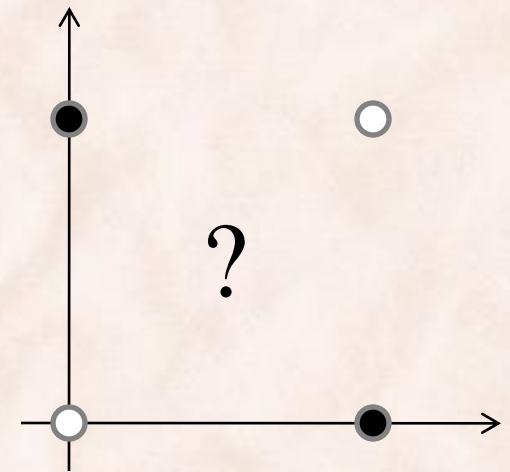
---



And



Or



Xor

$$\left. \begin{aligned} \mathbf{x} &= [1, x_1, x_2]^T \\ \mathbf{w} &= [w_0, w_1, w_2]^T \end{aligned} \right\} \Rightarrow \mathbf{w}^T \mathbf{x} = w_1 x_1 + w_2 x_2 + w_0$$

# Perceptron with Quadratic Kernel

---

- Discriminant function:

$$h(\mathbf{x}) = \sum_{j=1..N} \alpha_j y_j \varphi(\mathbf{x}_j)^T \varphi(\mathbf{x}) = \sum_{j=1..N} \alpha_j y_j K(\mathbf{x}_j, \mathbf{x})$$

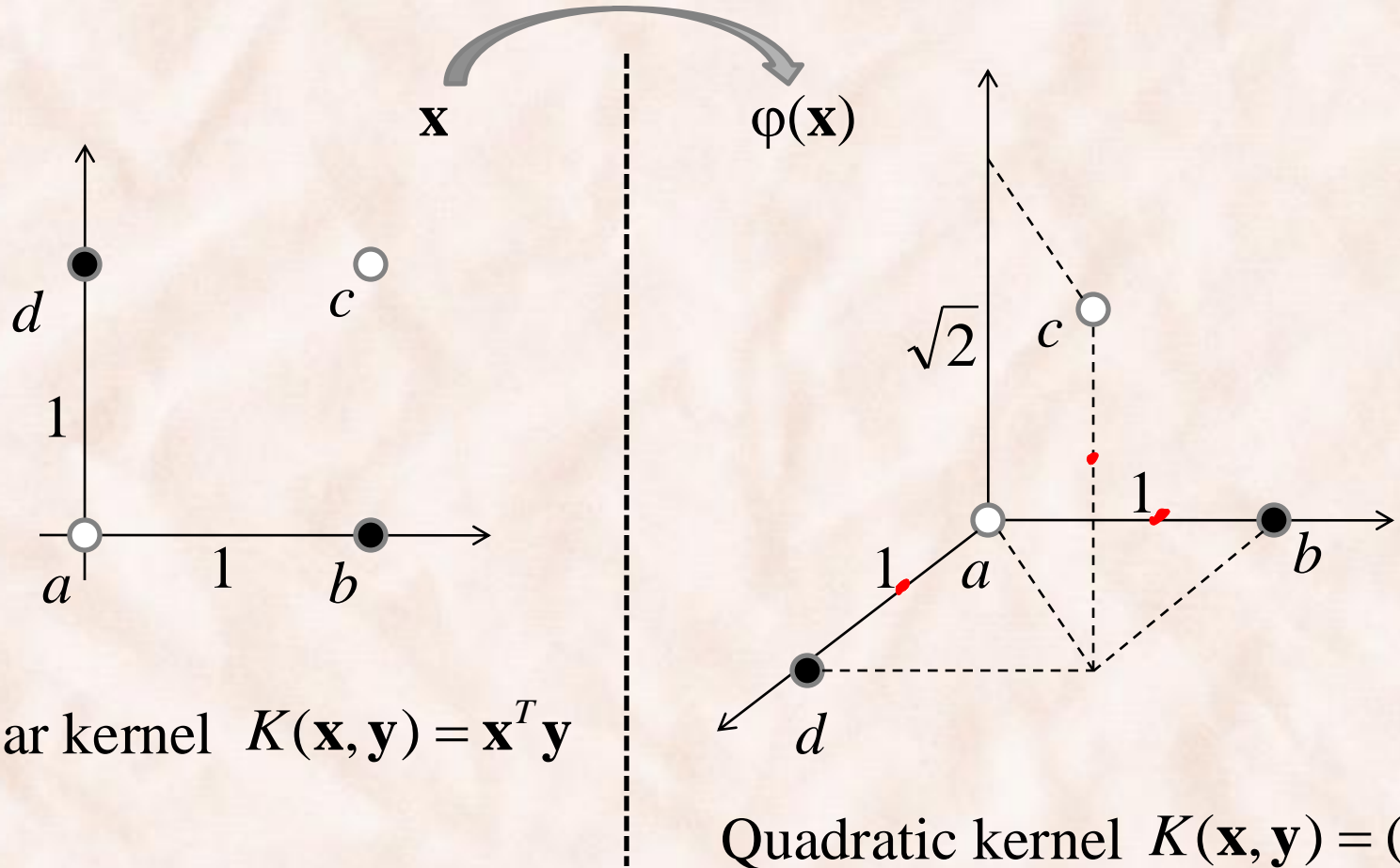
- Quadratic kernel:

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2 = (x_1 y_1 + x_2 y_2)^2$$

⇒ corresponding feature space  $\varphi(\mathbf{x}) = ?$

*conjunctions of two atomic features*

# Perceptron with Quadratic Kernel

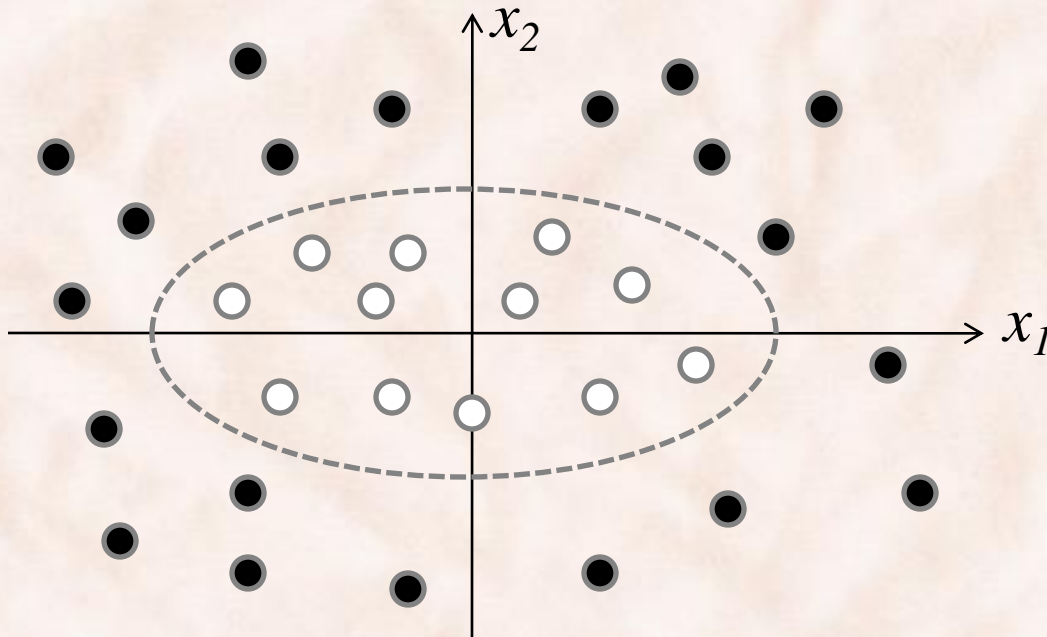


# Quadratic Kernels

- Circles, hyperbolas, and ellipses as separating surfaces:

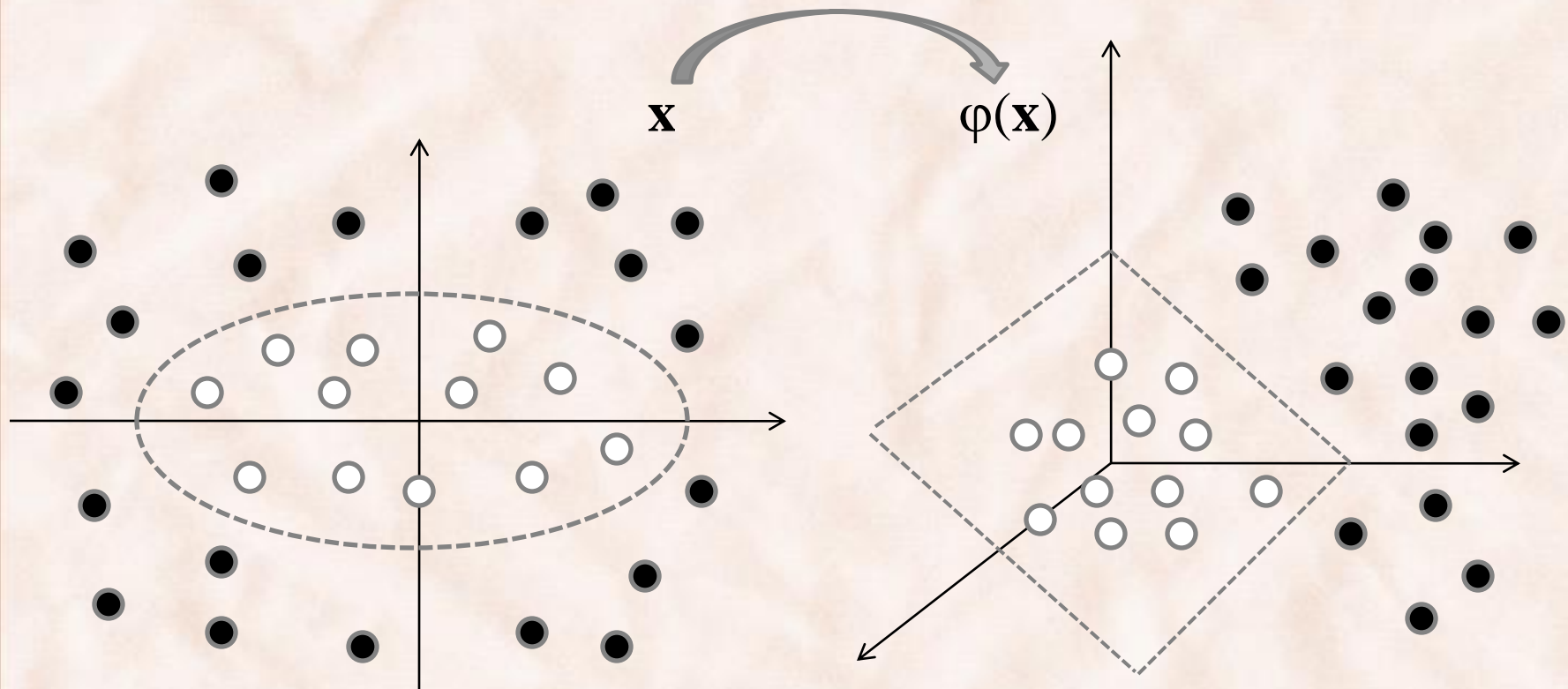
$$K(\mathbf{x}, \mathbf{y}) = (1 + \mathbf{x}^T \mathbf{y})^2 = \varphi(x)^T \varphi(y)$$

$$\varphi(x) = [1, \sqrt{2}x_1, \sqrt{2}x_2, x_1^2, \sqrt{2}x_1x_2, x_2^2]^T$$



# Quadratic Kernels

$$K(\mathbf{x}, \mathbf{y}) = (\mathbf{x}^T \mathbf{y})^2 = \varphi(\mathbf{x})^T \varphi(\mathbf{y})$$



# Explicit Features vs. Kernels

---

- Explicitly enumerating features can be prohibitive:
  - 1,000 basic features for  $\mathbf{x}^T \mathbf{y} \Rightarrow 500,500$  quadratic features for  $(\mathbf{x}^T \mathbf{y})^2$
  - Much worse for higher order features.
- Solution:
  - Do not compute the feature vectors, compute kernels instead (i.e. compute dot products between implicit feature vectors).
    - $(\mathbf{x}^T \mathbf{y})^2$  takes 1001 multiplications.
    - $\varphi(\mathbf{x})^T \varphi(\mathbf{y})$  in feature space takes 500,500 multiplications.



# Kernel Functions

---

- **Definition:**

A function  $k : X \times X \rightarrow \mathbb{R}$  is a kernel function if there exists a feature mapping  $\varphi : X \rightarrow \mathbb{R}^n$  such that:

$$k(\mathbf{x}, \mathbf{y}) = \varphi(\mathbf{x})^T \varphi(\mathbf{y})$$

- **Theorem:**

$k : X \times X \rightarrow \mathbb{R}$  is a valid kernel  $\Leftrightarrow$  the Gram matrix  $K$  whose elements are given by  $k(\mathbf{x}_n, \mathbf{x}_m)$  is *positive semidefinite* for all possible choices of the set  $\{\mathbf{x}_n\}$ .

# Kernel Examples

---

- **Linear kernel:**  $K(\mathbf{x}, \mathbf{y}) = \mathbf{x}^T \mathbf{y}$
- **Quadratic kernel:**  $K(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^2$ 
  - contains constant, linear terms and terms of order two ( $c > 0$ ).
- **Polynomial kernel:**  $K(\mathbf{x}, \mathbf{y}) = (c + \mathbf{x}^T \mathbf{y})^M$ 
  - contains all terms up to degree  $M$  ( $c > 0$ ).
- **Gaussian kernel:**  $K(\mathbf{x}, \mathbf{y}) = \exp(-\|\mathbf{x} - \mathbf{y}\|^2 / 2\sigma^2)$ 
  - Corresponding feature space has infinite dimensionality.
  - Prove using Taylor expansion of exponential.

also called  $r$  or  $\gamma$

$$\varphi(x) = e^{-\gamma x^2} [1, \sqrt{2\gamma}x, \sqrt{2\gamma}x^2, \dots]$$

# Techniques for Constructing Kernels

Given valid kernels  $k_1(\mathbf{x}, \mathbf{x}')$  and  $k_2(\mathbf{x}, \mathbf{x}')$ , the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \quad (6.13)$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}') \quad (6.14)$$

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.15)$$

$$k(\mathbf{x}, \mathbf{x}') = \exp(k_1(\mathbf{x}, \mathbf{x}')) \quad (6.16)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}') \quad (6.17)$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}') \quad (6.18)$$

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}')) \quad (6.19)$$

$$k(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{A} \mathbf{x}' \quad (6.20)$$

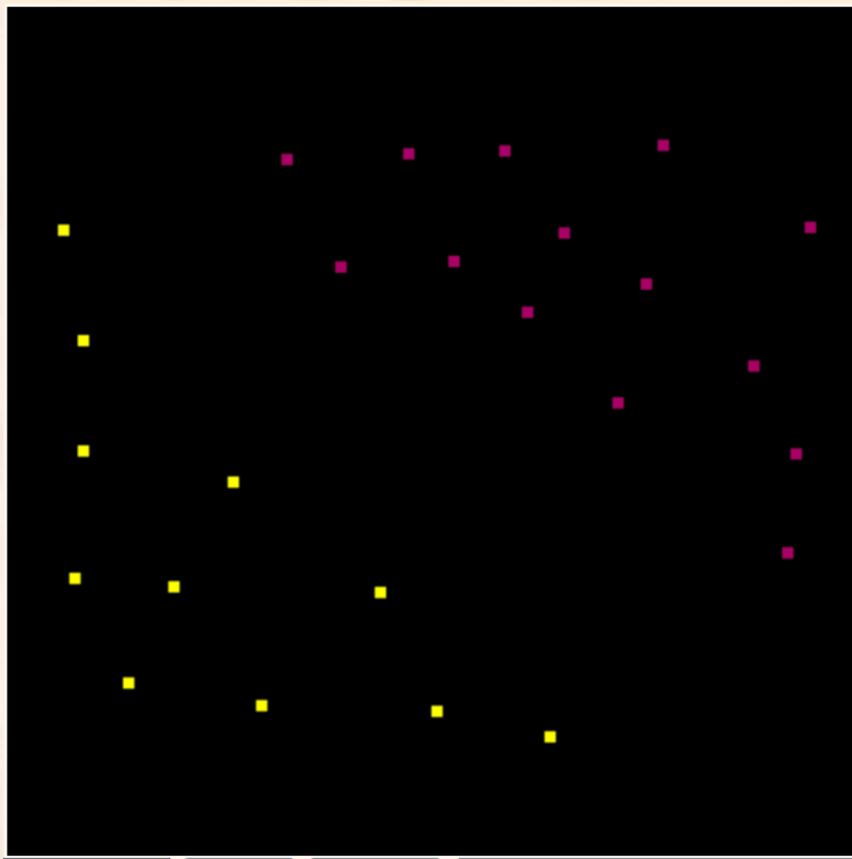
$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.21)$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b) \quad (6.22)$$

where  $c > 0$  is a constant,  $f(\cdot)$  is any function,  $q(\cdot)$  is a polynomial with nonnegative coefficients,  $\phi(\mathbf{x})$  is a function from  $\mathbf{x}$  to  $\mathbb{R}^M$ ,  $k_3(\cdot, \cdot)$  is a valid kernel in  $\mathbb{R}^M$ ,  $\mathbf{A}$  is a symmetric positive semidefinite matrix,  $\mathbf{x}_a$  and  $\mathbf{x}_b$  are variables (not necessarily disjoint) with  $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$ , and  $k_a$  and  $k_b$  are valid kernel functions over their respective spaces.

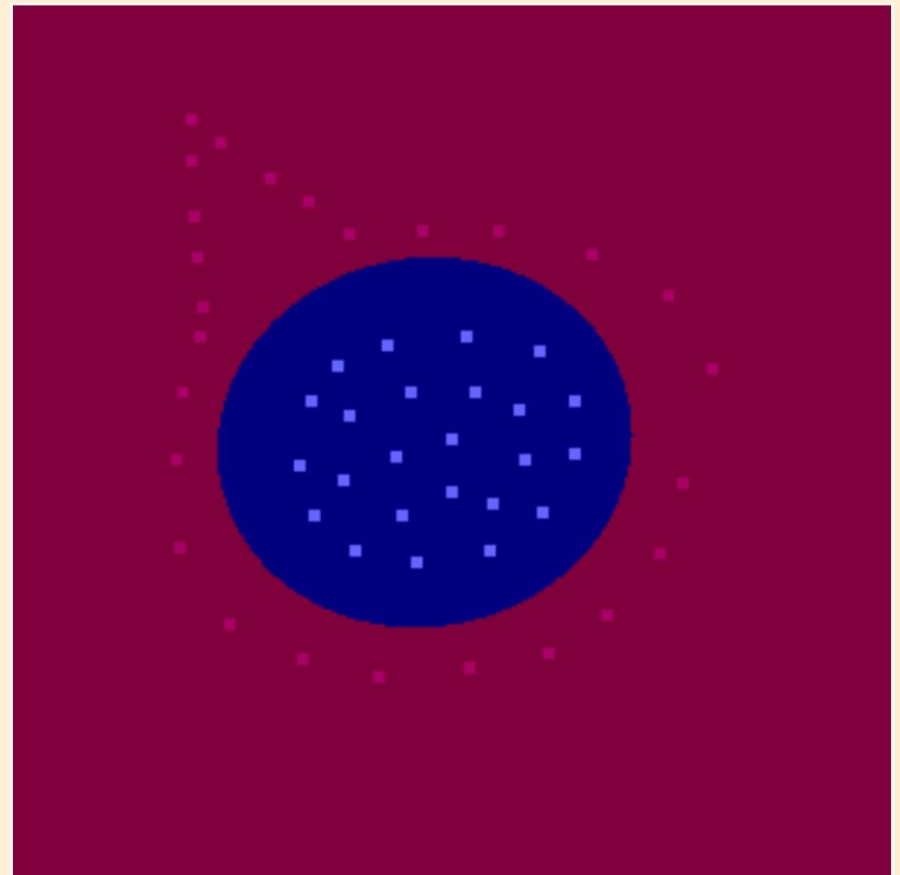
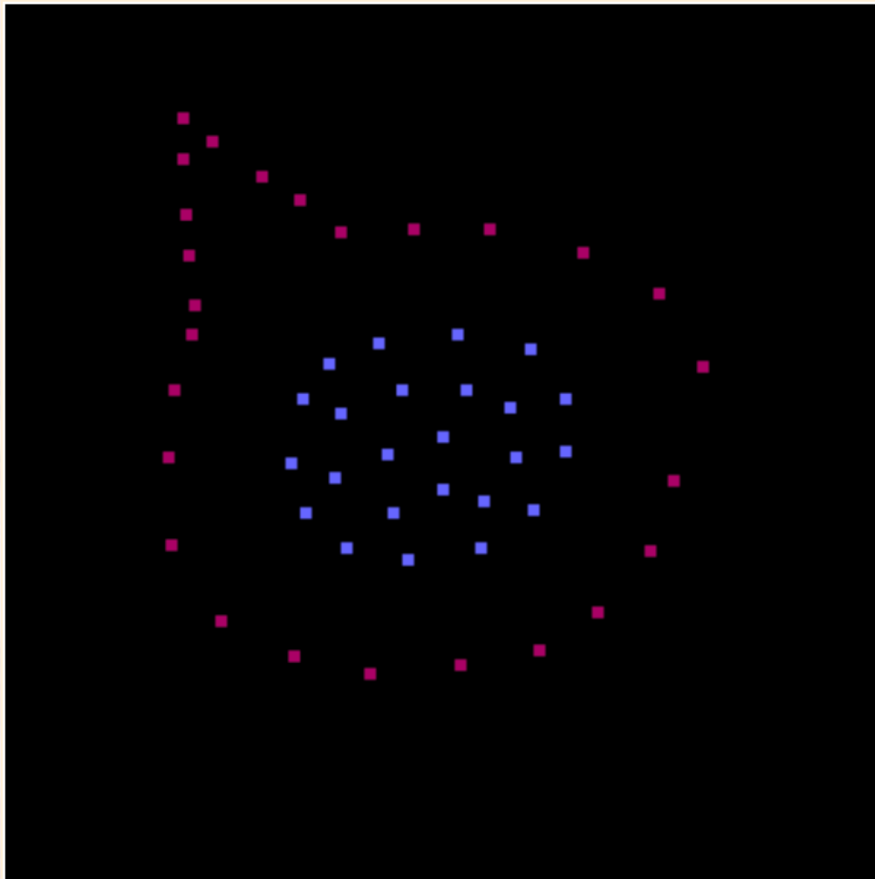
# Linear Kernel $\mathbf{x}^T \mathbf{y}$ , $C = 100$

---



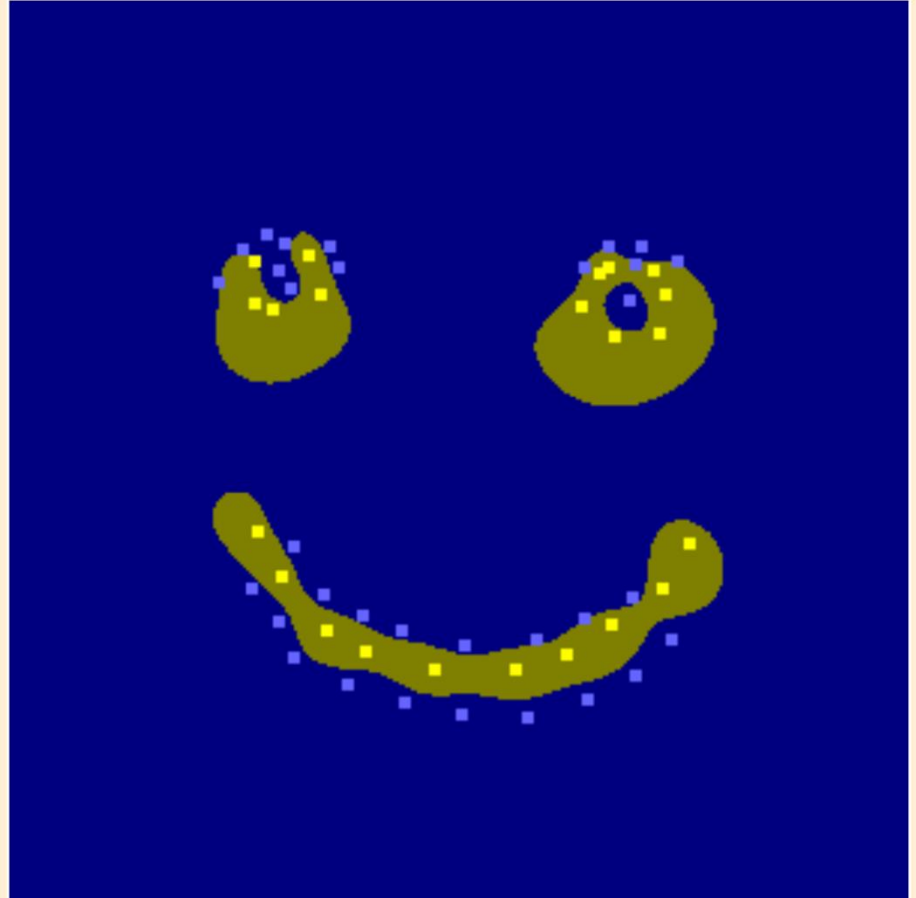
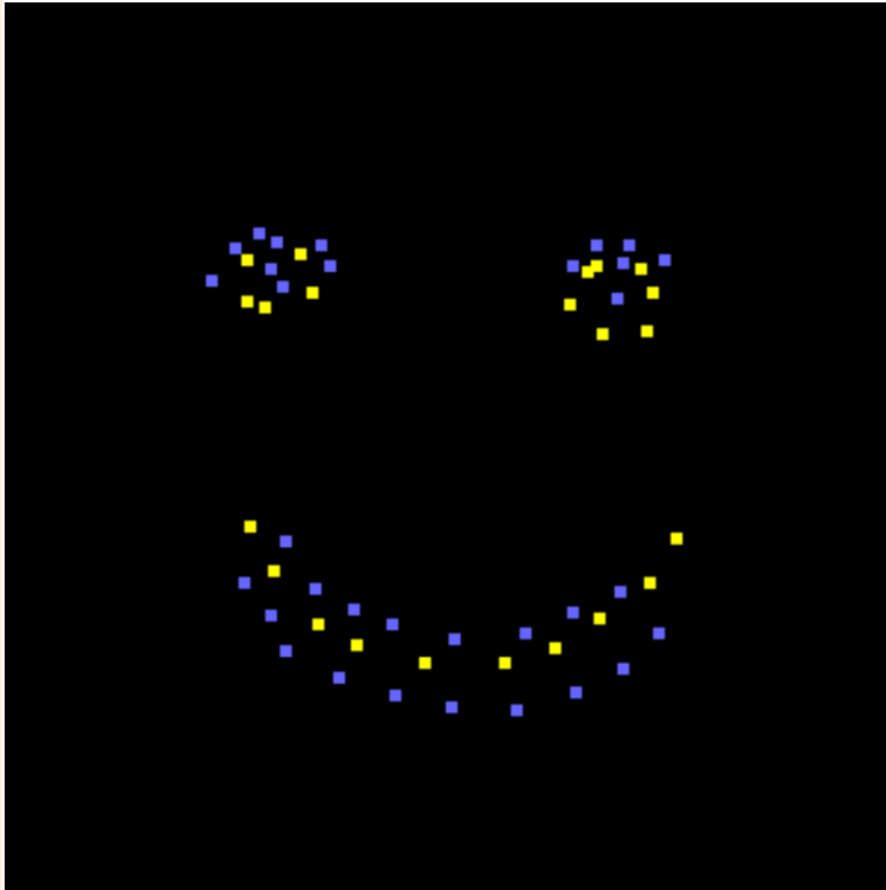
# Quadratic kernel $(1 + \mathbf{x}^T \mathbf{y})^2, C = 500$

---



Gaussian kernel with  $\frac{1}{2\sigma^2} = 250, C = 100$

---



# Kernels over Discrete Structures

---

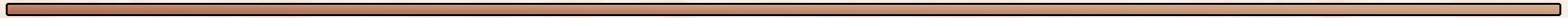
- **Subsequence Kernels** [Lodhi et al., JMLR 2002]:
  - $\Sigma$  is a finite alphabet (set of symbols).
  - $\mathbf{x}, \mathbf{y} \in \Sigma^*$  are two sequences of symbols with lengths  $|\mathbf{x}|$  and  $|\mathbf{y}|$
  - $k(\mathbf{x}, \mathbf{y})$  is defined as the number of common substrings of length  $n$ .
  - $k(\mathbf{x}, \mathbf{y})$  can be computed in  $O(n|\mathbf{x}||\mathbf{y}|)$  time complexity.
- **Tree Kernels** [Collins and Duffy, NIPS 2001]:
  - $T_1$  and  $T_2$  are two trees with  $N_1$  and  $N_2$  nodes respectively.
  - $k(T_1, T_2)$  is defined as the number of common subtrees.
  - $k(T_1, T_2)$  can be computed in  $O(N_1 N_2)$  time complexity.
  - in practice, time is linear in the size of the trees.

# Readings

---

1. [Required] Peter Flach's ML textbook, chapter 7:
  - Section 7.1 on the Least-squares method.
  - Section 7.2 on the Perceptron.
  - Section 7.5 on Kernels and the Kernel Perceptron.
2. [Optional] Bishop, chapter 6:
  - Section 6.1 on dual representations for linear regression models.
  - Section 6.2 on techniques for constructing new kernels.





# Linear Discriminant Functions: Multiple Classes ( $K > 2$ )

---

- 1) Train  $K$  or  $K-1$  *one-versus-the-rest* binary classifiers.
- 2) Train  $K(K-1)/2$  *one-versus-one* binary classifiers.

- 3) Train  $K$  linear functions:

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \boldsymbol{\varphi}(\mathbf{x}) + w_{k0}$$

- Decision:

$\mathbf{x} \in C_k$  if  $y_k(\mathbf{x}) > y_j(\mathbf{x})$ , for all  $j \neq k$ .

$\Rightarrow$  decision boundary between classes  $C_k$  and  $C_j$  is hyperplane defined

by  $y_k(\mathbf{x}) = y_j(\mathbf{x})$  i.e.  $(\mathbf{w}_k - \mathbf{w}_j)^T \boldsymbol{\varphi}(\mathbf{x}) + (w_{k0} - w_{j0}) = 0$

$\Rightarrow$  same geometrical properties as in binary case.

# Linear Discriminant Functions: Multiple Classes ( $K > 2$ )

---

4) More general ranking approach:

$$y(\mathbf{x}) = \arg \max_{t \in T} \mathbf{w}^T \varphi(\mathbf{x}, t) \quad \text{where} \quad T = \{c_1, c_2, \dots, c_K\}$$

- It subsumes the approach with  $K$  separate linear functions.
- Useful when  $T$  is very large (e.g. exponential in the size of input  $\mathbf{x}$ ), assuming inference can be done efficiently.

# The Perceptron Algorithm: K classes

---

1. **initialize** parameters  $\mathbf{w} = 0$
2. **for**  $i = 1 \dots n$
3.      $y_i = \arg \max_{t \in T} \mathbf{w}^T \varphi(\mathbf{x}_i, t)$
4.     **if**  $y_i \neq t_i$  **then**
5.          $\mathbf{w} = \mathbf{w} + \varphi(\mathbf{x}_i, t_i) - \varphi(\mathbf{x}_i, y_i)$

Repeat:

- a) until convergence.
- b) for a number of epochs E.

During testing:

$$t^* = \arg \max_{t \in T} \mathbf{w}^T f(\mathbf{x}, t)$$

# Averaged Perceptron: K classes

---

1. **initialize** parameters  $\mathbf{w} = 0$ ,  $\tau = 1$ ,  $\bar{\mathbf{w}} = 0$
2. **for**  $i = 1 \dots n$
3.      $y_i = \arg \max_{t \in T} \mathbf{w}^T \varphi(\mathbf{x}_i, t)$
4.     **if**  $y_i \neq t_i$  **then**
5.          $\mathbf{w} = \mathbf{w} + \varphi(\mathbf{x}_i, t_i) - \varphi(\mathbf{x}_i, y_i)$
6.          $\bar{\mathbf{w}} = \bar{\mathbf{w}} + \mathbf{w}$
7.          $\tau = \tau + 1$
8. **return**  $\bar{\mathbf{w}} / \tau$

Repeat:

- a) until convergence.
- b) for a number of epochs E.

During testing:  $t^* = \arg \max_{t \in T} \bar{\mathbf{w}}^T \varphi(\mathbf{x}, t)$

# The Perceptron Algorithm: K classes

---

1. **initialize** parameters  $\mathbf{w} = 0$
2. **for**  $i = 1 \dots n$
3.      $c_j = \arg \max_{t \in T} \mathbf{w}^T \varphi(\mathbf{x}_i, t)$
4.     **if**  $c_j \neq t_i$  **then**
5.          $\mathbf{w} = \mathbf{w} + \varphi(\mathbf{x}_i, t_i) - \varphi(\mathbf{x}_i, c_j)$

Repeat:

- a) until convergence.
- b) for a number of epochs E.

Loop invariant:  $\mathbf{w}$  is a weighted sum of training vectors:

$$\mathbf{w} = \hat{a} \sum_{i,j} a_{ij} (f(\mathbf{x}_i, t_i) - f(\mathbf{x}_i, c_j))$$

$$\Rightarrow \mathbf{w}^T f(\mathbf{x}, t) = \hat{a} \sum_{i,j} a_{ij} (f(\mathbf{x}_i, t_i)^T f(\mathbf{x}, t) - f(\mathbf{x}_i, c_j)^T f(\mathbf{x}, t))$$

# Kernel Perceptron: K classes

---

1. **define**  $f(\mathbf{x}, t) = \hat{a} \sum_{i,j} a_{ij} (f(\mathbf{x}_i, t_i)^T f(\mathbf{x}, t) - f(\mathbf{x}_i, c_j)^T f(\mathbf{x}, t))$
2. **initialize** dual parameters  $\alpha_{ij} = 0$
3. **for**  $i = 1 \dots n$
4.      $c_j = \arg \max_{t \in T} f(\mathbf{x}_i, t)$
5.     **if**  $y_i \neq t_i$  **then**
6.          $\alpha_{ij} = \alpha_{ij} + 1$

Repeat:

- a) until convergence.
- b) for a number of epochs E.

During testing:

$$t^* = \arg \max_{t \in T} f(\mathbf{x}, t)$$

# Kernel Perceptron: K classes

---

- Discriminant function:

$$\begin{aligned} f(\mathbf{x}, t) &= \hat{a}_{i,j} (f(\mathbf{x}_i, t_i)^T f(\mathbf{x}, t) - f(\mathbf{x}_i, c_j)^T f(\mathbf{x}, t)) \\ &= \hat{a}_{i,j} (K(\mathbf{x}_i, t_i, \mathbf{x}, t) - K(\mathbf{x}_i, c_j, \mathbf{x}, t)) \end{aligned}$$

where:

$$K(\mathbf{x}_i, t_i, \mathbf{x}, t) = \varphi^T(\mathbf{x}_i, t_i) \varphi(\mathbf{x}, t)$$

$$K(\mathbf{x}_i, y_i, \mathbf{x}, t) = \tilde{f}^T(\mathbf{x}_i, y_i) f(\mathbf{x}, t)$$



