

# Statistical Machine Learning

## Lecture 4: Classification

**Simone Schaub-Meyer & Marcus Rohrbach**

**Department of Computer Science**

**TU Darmstadt**

Summer Term 2025

# Organization - Homeworks

- Participating in the homework exercises is optional.
  - If you sign up be responsive and a team player.
  - Testate is only relevant for the corresponding homework.
  - No passing needed for participating in exam.
- Testate has to be done in person.
  - Part of the grade, i.e., same rules as for exams apply.
  - Only remote in exceptional cases (e.g., visa issues, illness) and with corresponding proof possible to do it online.
  - Sign up as group for a slot. Do not mix groups.
- Exercises are hard.
  - They also get you a bonus of up to 1.0 for the exam.
  - Partial bonus is possible.

# Classification

- What do you remember about classification?

# Today's Objectives

- Make you understand how to do build a discriminative classifier!
- Covered Topics:
  - Discriminant Functions
  - Multi-Class Classification
  - Fisher Discriminate Analysis
  - Perceptrons
  - Logistic Regression

# Outline

- 1. Discriminant Functions**
- 2. Fisher Discriminant Analysis**
- 3. Perceptron Algorithm**
- 4. Logistic Regression**
- 5. Wrap-Up**

# Outline

## 1. Discriminant Functions

## 2. Fisher Discriminant Analysis

## 3. Perceptron Algorithm

## 4. Logistic Regression

## 5. Wrap-Up

# Reminder of Bayesian Decision Theory

- We want to find the **a-posteriori probability** (posterior) of the class  $C_k$  given the observation (feature)  $x$

$$p(C_k | x) = \frac{p(x | C_k)p(C_k)}{p(x)} = \frac{p(x | C_k)p(C_k)}{\sum_j p(x | C_j)p(C_j)}$$

- $p(C_k | x)$  - class posterior
- $p(x | C_k)$  - class-conditional probability (likelihood)
- $p(C_k)$  - class prior
- $p(x)$  - normalization term

# Reminder of Bayesian Decision Theory

## ■ Decision rule

- Decide  $C_1$  if  $p(C_1 | x) > p(C_2 | x)$
- Using the definition of conditional distributions, equivalent to

$$p(x | C_1) p(C_1) > p(x | C_2) p(C_2) \equiv \frac{p(x | C_1)}{p(x | C_2)} > \frac{p(C_2)}{p(C_1)}$$

- A classifier obeying this rule is called a **Bayes optimal classifier**



# Reminder of Bayesian Decision Theory

## ■ Current approach

- $p(C_k | x) = p(x | C_k) p(C_k) / p(x)$  (Bayes' rule)
- Model and estimate the class-conditional density  $p(x | C_k)$  and the class prior  $p(C_k)$
- Compute posterior  $p(C_k | x)$
- Minimize the error probability by maximizing  $p(C_k | x)$

## ■ New approach

- Directly encode the *decision boundary*
- Without modeling the densities directly
- Still minimize the error probability

# Discriminant Functions

- Formulate classification using comparisons

- Discriminant functions

$$y_1(x), \dots, y_K(x)$$

- Classify  $x$  as class  $C_k$  iff

$$y_k(x) > y_j(x) \quad \forall j \neq k$$

- More formally, a **discriminant** maps a vector  $\mathbf{x}$  to one of the  $K$  available classes

# Discriminant Functions

- Example of discriminant functions from the Bayes classifier

$$y_k(x) = p(C_k | x)$$

$$y_k(x) = p(x | C_k) p(C_k)$$

$$y_k(x) = \log p(x | C_k) + \log p(C_k)$$

# Discriminant Functions

- Base case with 2 classes

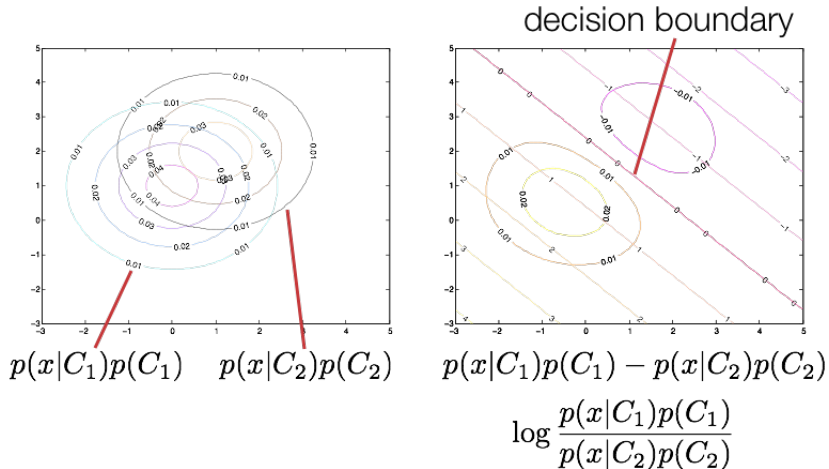
$$\begin{aligned}y_1(x) &> y_2(x) \\ y_1(x) - y_2(x) &> 0 \\ y(x) &> 0\end{aligned}$$

- Example from the Bayes classifier

$$\begin{aligned}y(x) &= p(C_1 | x) - p(C_2 | x) \\ y(x) &= \log \frac{p(x | C_1)}{p(x | C_2)} + \log \frac{p(C_1)}{p(C_2)}\end{aligned}$$

# Example - Bayes Classifier

- Base case with 2 classes and Gaussian class-conditionals



# Linear Discriminant Functions

- Base case with 2 classes

$y(\mathbf{x}) > 0$  decide class 1, otherwise class 2

- Simplest case: linear decision boundary
  - In **linear discriminants**, the decision surfaces are (hyper)planes
  - **Linear Discriminant Function**

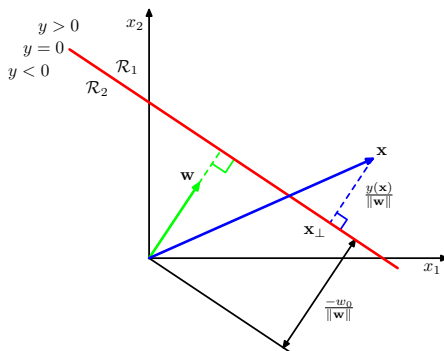
$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

- Where  $\mathbf{w}$  is the normal vector and  $w_0$  the offset

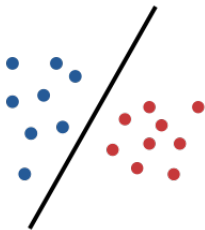
# Linear Discriminant Functions

## ■ Illustration of the 2D case

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0, \quad \mathbf{x} = \begin{bmatrix} x_1 & x_2 \end{bmatrix}^T$$



# Linear Discriminant Functions



Linearly separable



Not linearly separable



# Discriminant Functions

- Why might we want to use discriminant functions?



- We could easily fit the class-conditionals using Gaussians and use a Bayes classifier

# Discriminant Functions

- How about now? Do these points matter for making the decision between the two classes?



# Distribution-free Classifiers

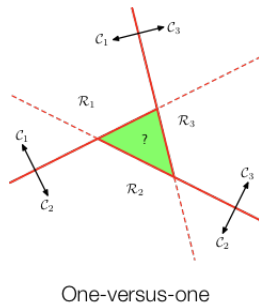
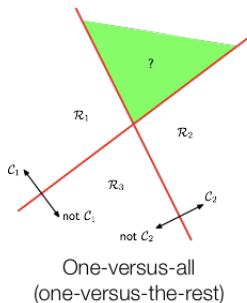
- We do not necessarily need to model all the details of the class-conditional distributions to come up with a good decision boundary. (The class-conditionals may have many intricacies that do not matter at the end of the day)



- If we can learn where to place the decision boundary directly, we can avoid some of the complexity
- It would be unwise to believe that such classifiers are inherently superior to probabilistic ones. We shall see why later...

# Multi-Class Case

- What if we constructed a multi-class classifier from several 2-class classifiers?



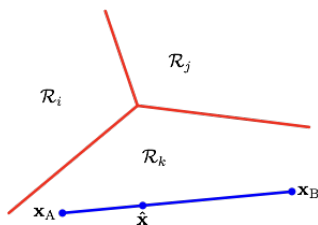
- If we base our decision rule on binary decisions, this may lead to ambiguities, where we can vote for several classes such as  $\mathcal{C}_1, \mathcal{C}_2$  respectively  $\mathcal{C}_1, \mathcal{C}_2, \mathcal{C}_3$

# Multi-Class Case - Better Solution

- Use a discriminant function to encode how strongly we believe in each class

$$y_1(x), \dots, y_K(x)$$

- **Decision rule:** Decide  $k$  if  $y_k(x) > y_j(x) \quad \forall j \neq k$



- If the discriminant functions are linear, the decision regions are **connected** and **convex**

# Outline

1. Discriminant Functions

**2. Fisher Discriminant Analysis**

3. Perceptron Algorithm

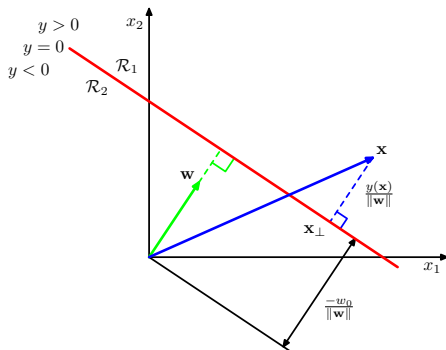
4. Logistic Regression

5. Wrap-Up

# Linear Discriminant Functions

## ■ Illustration of the 2D case

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0, \quad \mathbf{x} = \begin{bmatrix} x_1 & x_2 \end{bmatrix}^T$$



# First Attempt: Least Squares

- Try to achieve a certain value of the discriminative function

$$y(\mathbf{x}) = +1 \Leftrightarrow \mathbf{x} \in C_1$$

$$y(\mathbf{x}) = -1 \Leftrightarrow \mathbf{x} \in C_2$$

- Training data inputs:  $X = \{\mathbf{x}_1 \in \mathbb{R}^d, \dots, \mathbf{x}_n\}$
- Training data labels:  $Y = \{y_1 \in \{-1, +1\}, \dots, y_n\}$
- Linear Discriminant Function
  - Try to enforce  $\mathbf{x}_i^T \mathbf{w} + w_0 = y_i, \quad \forall i = 1, \dots, n$
  - There is one linear equation for each training data point/label pair



# First Attempt: Least Squares

## ■ Linear system of equations

$$\mathbf{x}_i^\top \mathbf{w} + w_0 = y_i, \quad \forall i = 1, \dots, n$$

■ Define  $\hat{\mathbf{x}}_i = [\mathbf{x}_i \quad 1]^\top \in \mathbb{R}^{(d+1) \times 1}$ ,  $\hat{\mathbf{w}} = [\mathbf{w} \quad w_0]^\top \in \mathbb{R}^{(d+1) \times 1}$

■ Rewrite the equation system

$$\hat{\mathbf{x}}_i^\top \hat{\mathbf{w}} = y_i, \quad \forall i = 1, \dots, n$$

■ In **matrix-vector notation** we have

$$\hat{\mathbf{X}}^\top \hat{\mathbf{w}} = \mathbf{y}$$

■ With  $\hat{\mathbf{X}} = [\hat{\mathbf{x}}_1, \dots, \hat{\mathbf{x}}_n] \in \mathbb{R}^{(d+1) \times n}$  and  $\mathbf{y} = [y_1, \dots, y_n]^\top$

# First Attempt: Least Squares

$$\hat{\mathbf{X}}^T \hat{\mathbf{w}} = \mathbf{y}$$

- An overdetermined system of equations
- There are  $n$  equations and  $d + 1$  unknowns

# First Attempt: Least Squares

- Look for the least squares solution

$$\hat{\mathbf{w}}^* = \arg \min_{\hat{\mathbf{w}}} \left\| \hat{\mathbf{X}}^T \hat{\mathbf{w}} - \mathbf{y} \right\|^2$$

How to solve this problem?

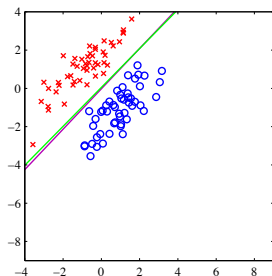
$$\begin{aligned} \hat{\mathbf{w}}^* &= \arg \min_{\hat{\mathbf{w}}} \left( \hat{\mathbf{X}}^T \hat{\mathbf{w}} - \mathbf{y} \right)^T \left( \hat{\mathbf{X}}^T \hat{\mathbf{w}} - \mathbf{y} \right) \\ &= \arg \min_{\hat{\mathbf{w}}} \hat{\mathbf{w}}^T \hat{\mathbf{X}} \hat{\mathbf{X}}^T \hat{\mathbf{w}} - 2 \mathbf{y}^T \hat{\mathbf{X}}^T \hat{\mathbf{w}} + \mathbf{y}^T \mathbf{y} \end{aligned}$$

$$\nabla_{\hat{\mathbf{w}}} \left( \hat{\mathbf{w}}^T \hat{\mathbf{X}} \hat{\mathbf{X}}^T \hat{\mathbf{w}} - 2 \mathbf{y}^T \hat{\mathbf{X}}^T \hat{\mathbf{w}} + \mathbf{y}^T \mathbf{y} \right) = 0$$

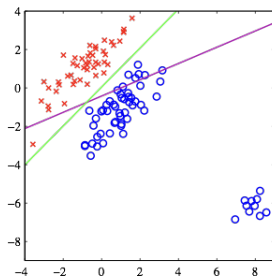
$$\implies \hat{\mathbf{w}} = \underbrace{\left( \hat{\mathbf{X}} \hat{\mathbf{X}}^T \right)^{-1}}_{\text{pseudo-inverse}} \hat{\mathbf{X}} \mathbf{y}$$

# First Attempt: Least Squares

- Problem: Least-squares is very sensitive to outliers



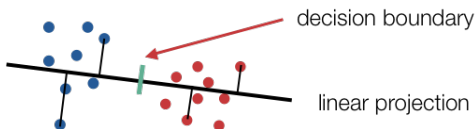
Without outliers  
least-squares discriminant  
works



With outliers least-squares  
discriminant breaks down

# Fisher's Linear Discriminant

- Take a different view on linear classification
- Find a **linear projection** of our data and classify the projected values



- The same thing as a linear discriminant function
  - Projection:  $y = \mathbf{w}^T \mathbf{x}$
  - Checking against a threshold:  $\mathbf{w}^T \mathbf{x} \geq -w_0$  or  $\mathbf{w}^T \mathbf{x} + w_0 \geq 0$

# Fisher's Linear Discriminant

- What is a **good projection  $\mathbf{w}$** ?
  - Idea: **Maximize the “distance” between the two classes to allow for a good separation**
- First attempt: **Maximize the distance between the class means**

$$\mathbf{m}_1 = \frac{1}{|C_1|} \sum_{i \in C_1} \mathbf{x}_i \quad \mathbf{m}_2 = \frac{1}{|C_2|} \sum_{i \in C_2} \mathbf{x}_i$$

- Projection of the means on the 1D line of real numbers

$$m_1 = \mathbf{w}^T \mathbf{m}_1 \quad m_2 = \mathbf{w}^T \mathbf{m}_2$$

- Maximize squared distance between means

$$\max (m_1 - m_2)^2$$

# Fisher's Linear Discriminant

- Maximize squared distance between means

$$\mathbf{w}^* = \arg \max_{\mathbf{w}} (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2$$

Can this maximization problem be solved?

- Obvious problem: Grows unboundedly with the norm of  $\mathbf{w}$
- Obvious solution: Fix the norm of  $\mathbf{w}$

$$\begin{aligned} \max_{\mathbf{w}} \quad & (\mathbf{w}^T \mathbf{m}_1 - \mathbf{w}^T \mathbf{m}_2)^2 \\ \text{s.t.} \quad & \|\mathbf{w}\|^2 = 1 \end{aligned}$$

- Constrained optimization problem!

# Fisher's Linear Discriminant

$$\begin{aligned} \max_{\mathbf{w}} \quad & (\mathbf{w}^\top \mathbf{m}_1 - \mathbf{w}^\top \mathbf{m}_2)^2 \\ \text{s.t.} \quad & \|\mathbf{w}\|^2 = 1 \end{aligned}$$

- Necessary conditions

$$\begin{aligned} \nabla_{\mathbf{x}} f(\mathbf{x}) + \lambda \nabla_{\mathbf{x}} g(\mathbf{x}) &= 0 \\ 2(\mathbf{w}^\top \mathbf{m}_1 - \mathbf{w}^\top \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2) + 2\lambda \mathbf{w} &= 0 \end{aligned}$$

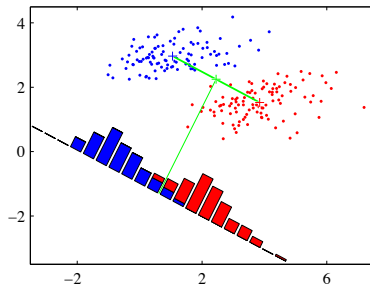
- It follows that

$$\mathbf{w} = \frac{\mathbf{m}_1 - \mathbf{m}_2}{\|\mathbf{m}_1 - \mathbf{m}_2\|}$$



# Fisher's Linear Discriminant

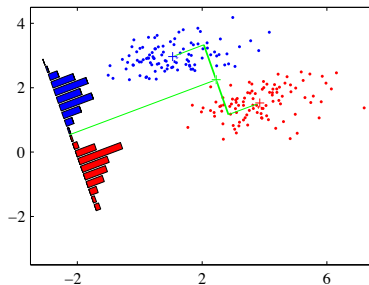
■ Here's what we get



■ Obvious problem: **large class overlap**

# Fisher's Linear Discriminant

- Here's what we **could** get



- Much better separation between classes
- How do we get this?
  - Idea: **Separate the means as far as possible while minimizing the variance of each class**

# Fisher's Linear Discriminant

- Second (and final) attempt:
  - Define **within-class variances**:

$$s_1^2 = \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 \quad s_2^2 = \sum_{n \in C_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2$$

where  $m_1 = \mathbf{w}^T \mathbf{m}_1$  and  $m_2 = \mathbf{w}^T \mathbf{m}_2$

- **Fisher criterion**

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

# Fisher's Linear Discriminant

## ■ Fisher criterion

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

## ■ Rewrite the numerator

$$\begin{aligned}
 (m_1 - m_2)^2 &= (\mathbf{w}^\top \mathbf{m}_1 - \mathbf{w}^\top \mathbf{m}_2)^2 \\
 &= (\mathbf{w}^\top (\mathbf{m}_1 - \mathbf{m}_2))^2 \\
 &= \mathbf{w}^\top \underbrace{(\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top}_{\substack{=: \mathbf{S}_B \\ \text{between-class covariance}}} \mathbf{w}
 \end{aligned}$$

# Fisher's Linear Discriminant

## ■ Fisher criterion

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2}$$

## ■ Rewrite the denominator

$$\begin{aligned}
 s_1^2 + s_2^2 &= \sum_{n \in C_1} (\mathbf{w}^T \mathbf{x}_n - m_1)^2 + \sum_{n \in C_2} (\mathbf{w}^T \mathbf{x}_n - m_2)^2 \\
 &= \sum_{n \in C_1} (\mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_1))^2 + \sum_{n \in C_2} (\mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_2))^2 \\
 &= \sum_{n \in C_1} \mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_1) (\mathbf{x}_n - \mathbf{m}_1)^T \mathbf{w} + \sum_{n \in C_2} \mathbf{w}^T (\mathbf{x}_n - \mathbf{m}_2) (\mathbf{x}_n - \mathbf{m}_2)^T \mathbf{w} \\
 &= \mathbf{w}^T \underbrace{\left[ \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1) (\mathbf{x}_n - \mathbf{m}_1)^T + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2) (\mathbf{x}_n - \mathbf{m}_2)^T \right]}_{\substack{=: \mathbf{S}_W \\ \text{within-class covariance}}} \mathbf{w}
 \end{aligned}$$

# Fisher's Linear Discriminant

## ■ Fisher criterion

$$\max_{\mathbf{w}} J(\mathbf{w}) = \frac{(m_1 - m_2)^2}{s_1^2 + s_2^2} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

- Differentiating w.r.t.  $\mathbf{w}$  and setting to 0 we have

$$(\mathbf{w}^T \mathbf{S}_B \mathbf{w}) \mathbf{S}_W \mathbf{w} = (\mathbf{w}^T \mathbf{S}_W \mathbf{w}) \mathbf{S}_B \mathbf{w}$$

What can we say about  $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$  and  $(\mathbf{w}^T \mathbf{S}_W \mathbf{w})$ ?

- Since  $(\mathbf{w}^T \mathbf{S}_B \mathbf{w})$  and  $(\mathbf{w}^T \mathbf{S}_W \mathbf{w})$  are scalars, we have that

$$\mathbf{S}_W \mathbf{w} \parallel \mathbf{S}_B \mathbf{w}$$

where  $\parallel$  means collinearity.

# Fisher's Linear Discriminant

- Also, we know that

$$\mathbf{S}_B \mathbf{w} = (\mathbf{m}_1 - \mathbf{m}_2)(\mathbf{m}_1 - \mathbf{m}_2)^\top \mathbf{w} \implies \mathbf{S}_B \mathbf{w} \parallel (\mathbf{m}_1 - \mathbf{m}_2)$$

- Hence, we have

$$\begin{aligned} \mathbf{S}_W \mathbf{w} &\parallel (\mathbf{m}_1 - \mathbf{m}_2) \\ \mathbf{w} &\parallel \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2) \end{aligned}$$

- Fisher's Linear Discriminant

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

# Fisher's Linear Discriminant

$$\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_1 - \mathbf{m}_2)$$

- The Fisher linear discriminant only gives us a **projection**
  - We still need to find the threshold
  - E.g., use Bayes classifier with Gaussian class-conditionals
- **Bayes optimality**
  - Fisher's linear discriminant is Bayes optimal if the class-conditional distributions are equal, with diagonal covariance
- Essentially equivalent to **Linear Discriminant Analysis** (LDA)



# Fisher's Linear Discriminant

- We won't go through this here, but Fisher's linear discriminant can be shown to be equivalent to a certain case of a least-squares linear classifier (see Bishop 4.1.5)
- **Problem with this method:** it is still very sensitive to noise!
- By The Way: This method is a **true classic** (it dates back to 1936)
  - Fisher, R.A., *The Use of Multiple Measurements in Taxonomic Problems*. Annals of Eugenics, 7: 179-188 (1936)

# Outline

1. Discriminant Functions
2. Fisher Discriminant Analysis
- 3. Perceptron Algorithm**
4. Logistic Regression
5. Wrap-Up

# New Strategy

- If our classes are **linearly separable**, we want to make sure that we find a **separating (hyper)plane**



- First such algorithm we will see
  - The perceptron algorithm [Rosenblatt, 1962]



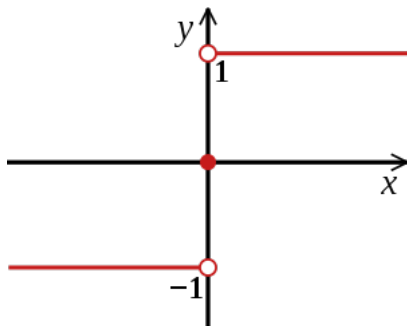
Rosenblatt [1928-1971]

# Perceptron Algorithm

## ■ Perceptron discriminant function

$$y(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} + b)$$

■ where  $\text{sign}(x) = \{+1, x > 0; 0, x = 0; -1, x < 0\}$

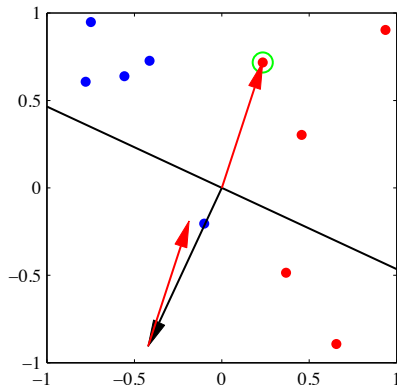


# Perceptron Algorithm

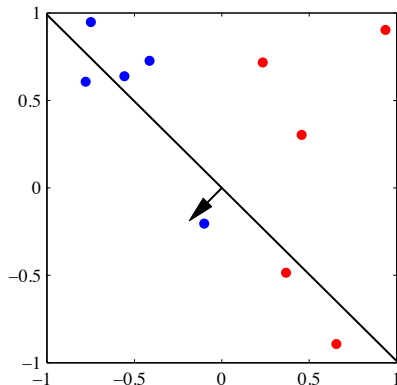
## ■ Perceptron Algorithm

- Initialize the weight vector  $\mathbf{w}$  and bias  $b$
- For all pairs of data points  $(\mathbf{x}_i, y_i)$ , where  $y_i \in \{-1, +1\}$ , do
  - If  $\mathbf{x}_i$  is correctly classified, i.e.,  $y(\mathbf{x}_i) = y_i$ , do nothing
  - Else if  $y_i = 1$  update the parameters with
$$\mathbf{w} \leftarrow \mathbf{w} + \mathbf{x}_i, \quad b \leftarrow b + 1$$
  - Else if  $y_i = -1$  update the parameters with
$$\mathbf{w} \leftarrow \mathbf{w} - \mathbf{x}_i, \quad b \leftarrow b - 1$$
- Repeat until convergence

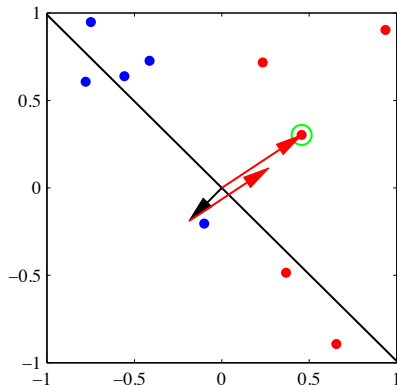
# Perceptron Algorithm - Intuition



# Perceptron Algorithm - Intuition

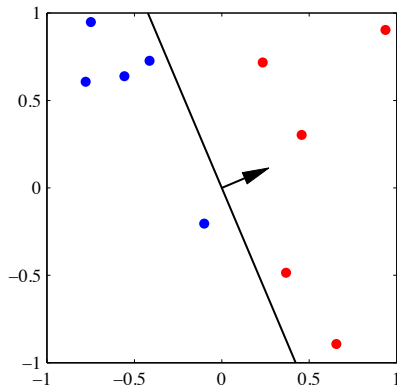


# Perceptron Algorithm - Intuition





# Perceptron Algorithm - Intuition



# Perceptron Algorithm

- Why does this algorithm work?
- We have an optimization problem

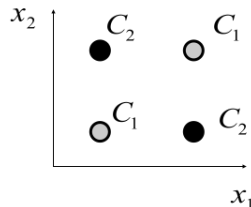
$$\begin{aligned}\max_{\mathbf{w}} J(\mathbf{w}) &= |\{\mathbf{x} \in C_1 : \mathbf{w}^\top \mathbf{x} > 0\}| + \dots \\ &= \sum_{\mathbf{x} \in C_1 : \mathbf{w}^\top \mathbf{x} < 0} \mathbf{w}^\top \mathbf{x} + \dots\end{aligned}$$

- And also a gradient method

$$\frac{\partial J}{\partial \mathbf{w}} = \sum_{\mathbf{x} \in C_1 : \mathbf{w}^\top \mathbf{x} < 0} \mathbf{x} + \dots$$

# But is the Perceptron Algorithm useful?

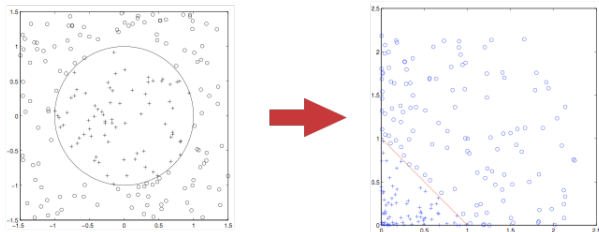
- How often is data linearly separable?
- A **simple failure** example is the XOR function



- History: Minsky & Papert [1969] criticized the perceptron for not being able to handle this case, which halted research on this and related techniques for decades

# Other Feature Spaces

- It took a long time until people had realized that there is a simple way out
- Key idea: Transform the input data nonlinearly so that the problem becomes linearly separable!



- There is an important message to get out from this
  - Create features instead of learning from raw data
  - Neural networks and kernel methods do it *automatically* for you

# Outline

1. Discriminant Functions
2. Fisher Discriminant Analysis
3. Perceptron Algorithm
- 4. Logistic Regression**
5. Wrap-Up

# Generative vs. Discriminative

- There are two different views to solve the classification problem
- **Generative modelling**
  - We model the class-conditional distributions  $p(x | C_2)$  and  $p(x | C_1)$
  - We classify by computing the class posterior using Bayes' rule
  - E.g.: Naive Bayes
- **Discriminative modelling**
  - We model the class-posterior directly, e.g.  $p(C_1 | x)$
  - Consequence: We only care about getting the classification right, and not whether we fit the class-conditional well
  - E.g.: Logistic Regression

# Probabilistic Discriminative Models

- For now, we will write the class posterior using Bayes' rule

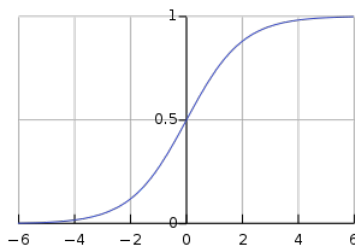
$$\begin{aligned} p(C_1 | \mathbf{x}) &= \frac{p(\mathbf{x} | C_1) p(C_1)}{p(\mathbf{x})} = \frac{p(\mathbf{x} | C_1) p(C_1)}{\sum_i p(\mathbf{x}, C_i)} \\ &= \frac{p(\mathbf{x} | C_1) p(C_1)}{\sum_i p(\mathbf{x} | C_i) p(C_i)} \\ &= \frac{p(\mathbf{x} | C_1) p(C_1)}{p(\mathbf{x} | C_1) p(C_1) + p(\mathbf{x} | C_2) p(C_2)} \\ &= \frac{1}{1 + p(\mathbf{x} | C_2) p(C_2) / (p(\mathbf{x} | C_1) p(C_1))} \\ &= \frac{1}{1 + \exp(-a)} = \sigma(a) \rightarrow \text{logistic sigmoid function} \end{aligned}$$

$$\text{with } a = \log \frac{p(\mathbf{x}|C_1)p(C_1)}{p(\mathbf{x}|C_2)p(C_2)}$$

# Sigmoid

## ■ Logistic / Sigmoid function

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$



[Wikipedia]

## ■ Sigmoid: 'S-shaped'

## ■ Squashes real numbers into the $[0, 1]$ interval



# Probabilistic Discriminative Models

- Class posterior

$$p(C_1 | \mathbf{x}) = \sigma(a) \quad \text{with} \quad a = \log \frac{p(\mathbf{x} | C_1) p(C_1)}{p(\mathbf{x} | C_2) p(C_2)}$$

- Logistic regression

- Assume that  $a$  is given by a linear discriminant function

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + w_0)$$

- Find  $\mathbf{w}$  and  $w_0$  so that the class-posterior is modeled best
  - When is this an **appropriate assumption**?
    - When the class conditionals are Gaussians with equal covariance
    - But also for a number of other distributions
    - Some independence of the form of the class-conditionals

# Logistic Regression

- Model the class posterior as

$$p(C_1 | \mathbf{x}) = \sigma(\mathbf{w}^\top \mathbf{x} + w_0)$$

- Maximize the likelihood

- Data (as always) is i.i.d. and define  $y_i = \begin{cases} 0 & \mathbf{x}_i \text{ belongs to } C_1 \\ 1 & \mathbf{x}_i \text{ belongs to } C_2 \end{cases}$

$$\begin{aligned} p(Y | X; \mathbf{w}, w_0) &= \prod_{i=1}^N p(y_i | \mathbf{x}_i; \mathbf{w}, w_0) \\ &= \prod_{i=1}^N p(C_1 | \mathbf{x}_i; \mathbf{w}, w_0)^{1-y_i} p(C_2 | \mathbf{x}_i; \mathbf{w}, w_0)^{y_i} \\ &= \prod_{i=1}^N \sigma(\mathbf{w}^\top \mathbf{x}_i + w_0)^{1-y_i} (1 - \sigma(\mathbf{w}^\top \mathbf{x}_i + w_0))^{y_i} \end{aligned}$$

# Logistic Regression

- We won't do the derivation here (see Bishop 4.3), but basically you can apply the logarithm to  $p(Y | X; \mathbf{w}, w_0)$  and do gradient descent
- Later, we will turn to a very different interpretation of this:
  - [Logistic regression as a neural network](#)

# Outline

1. Discriminant Functions
2. Fisher Discriminant Analysis
3. Perceptron Algorithm
4. Logistic Regression
- 5. Wrap-Up**

## 5. Wrap-Up

You know now:

- What a Bayesian Optimal Classifier is
- What a discriminant function is
- How to formalize (with intuition and mathematically) the classification problem as linearly-separable
- How to compute the least squares solution for classification and why it fails
- What Fisher's Linear Discriminant is and how it differs from least-squares
- What the perceptron is, why it fails in the XOR problem and how to overcome it with feature spaces
- The difference between Generative and Discriminative modelling
- What logistic regression is

# Self-Test Questions

- How do we get from Bayesian optimal decisions to discriminant functions?
- How to derive a discriminant function from a probability distribution?
- How to deal with more than two classes?
- What does linearly-separable mean?
- What is Fisher discriminant analysis? How does it relate to regression?
- Is Fisher's linear discriminant Bayes optimal?
- What are perceptrons? How can we train them?
- What is logistic regression? How to derive the parameter update rule?

# Homework

- To get a deeper understanding towards today's topics
  - Bishop, Chapter 4
  - The Elements of Statistical Learning, Chapter 3