

Machine Learning

Linear Discriminant Functions - Lecture V

Course Outline

Basic Concepts

- Parametric Method,
- Bayesian Learning and Nonparametrics Methods
- Clustering and Mixture of Gaussians

Classification Approaches

- Linear Discriminants
- Ensemble Methods and Boosting
- Randomized Trees, Forest

Deep Learning

- Foundations
- Optimization

Reinforcement Learning

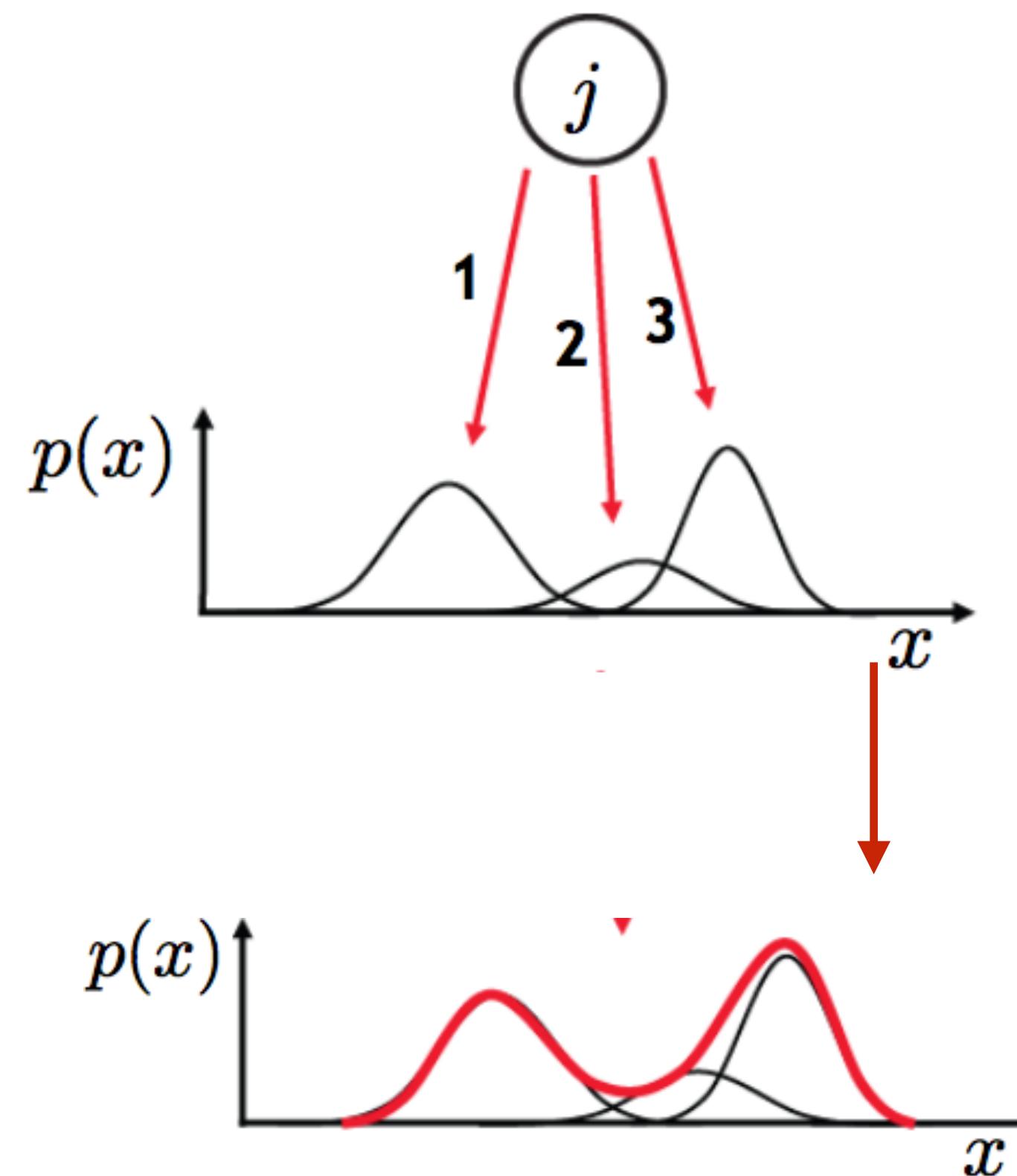
- Classical Reinforcement Learning
- Deep Reinforcement Learning

Videos for This Lecture

- Repetition Video (Part 0)
- Linear discriminant function (Part 1)
- Least-squares classification (Part 2)
- Generalized linear models (Part 3)

Recap: Mixture of Gaussians (MoG)

Generative Model



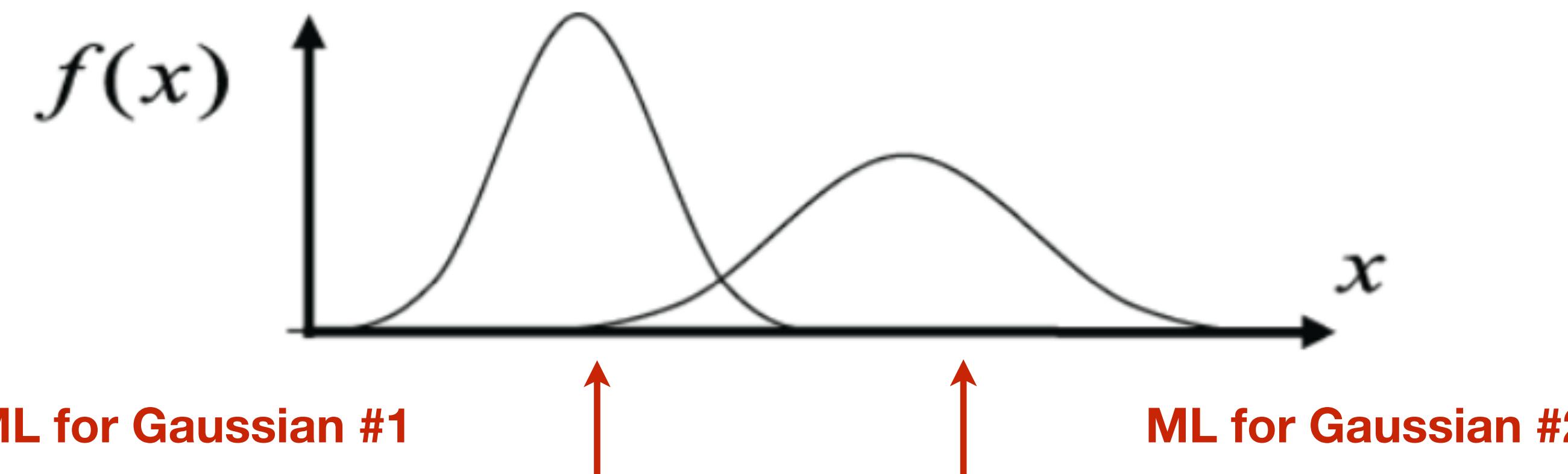
$$p(j) = \pi_j \quad \text{WEIGHT of mixture component}$$

$$p(x|\theta_j) \quad \text{Mixture component}$$

$$p(x|\theta) = \sum_{j=1}^M p(x|\theta_j)p(j) \quad \text{Mixture density}$$

Recap: Mixture of Gaussians

Assume we knew the values of the hidden variable...



assumed known \longrightarrow $1\ 1\ 1\ 1 \quad 2\ 2\ 2\ 2 \quad j$

$$h(j=1|x_n) = \begin{matrix} 1 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 \end{matrix}$$

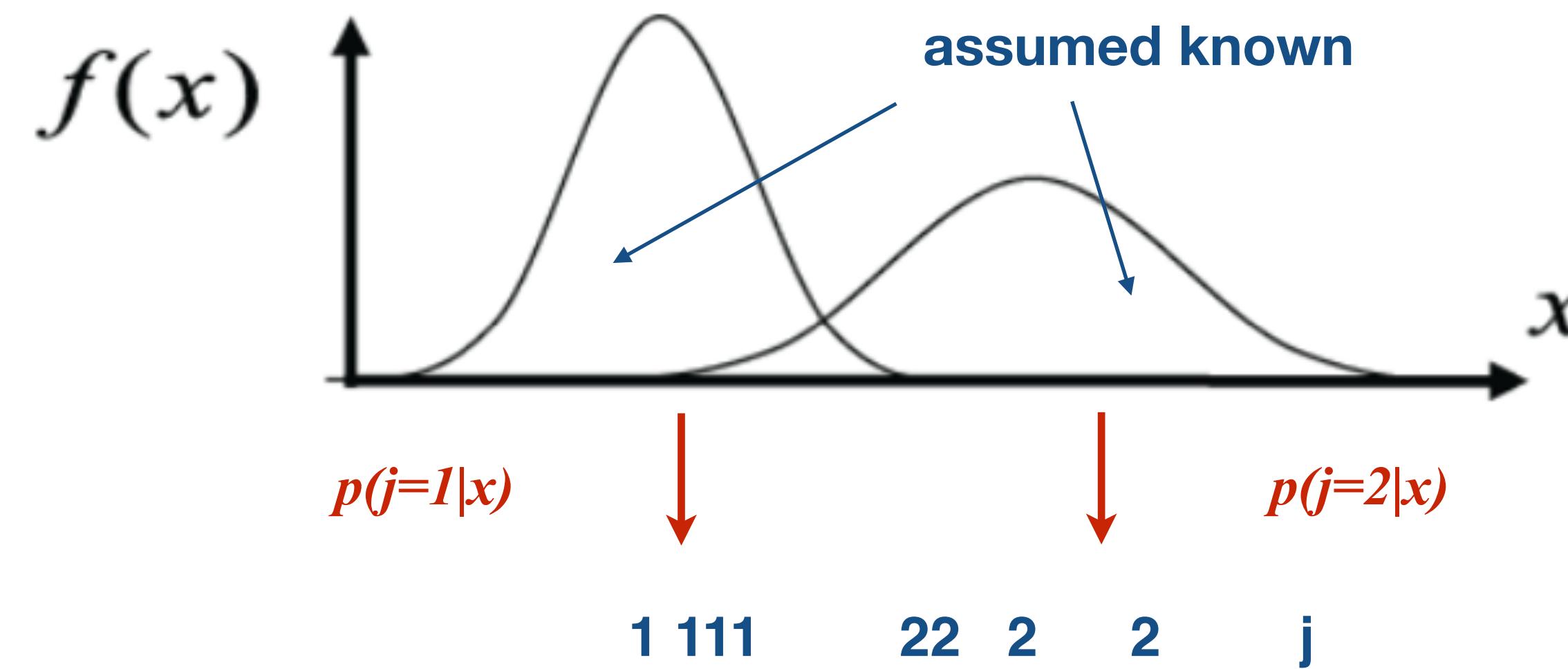
$$h(j=2|x_n) = \begin{matrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 1 & 1 \end{matrix}$$

$$\mu_1 = \frac{\sum_{n=1}^N h(j=1|x_n)x_n}{\sum_{i=1}^N h(j=1|x_n)}$$

$$\mu_2 = \frac{\sum_{n=1}^N h(j=2|x_n)x_n}{\sum_{i=1}^N h(j=2|x_n)}$$

Recap: Estimating MoGs-Iterative Strategy

Assume we knew the values of the hidden variable...



Bayes decision rule:

Decide $j = 1$ if $p(j = 1|x_n) > p(j = 2|x_n)$

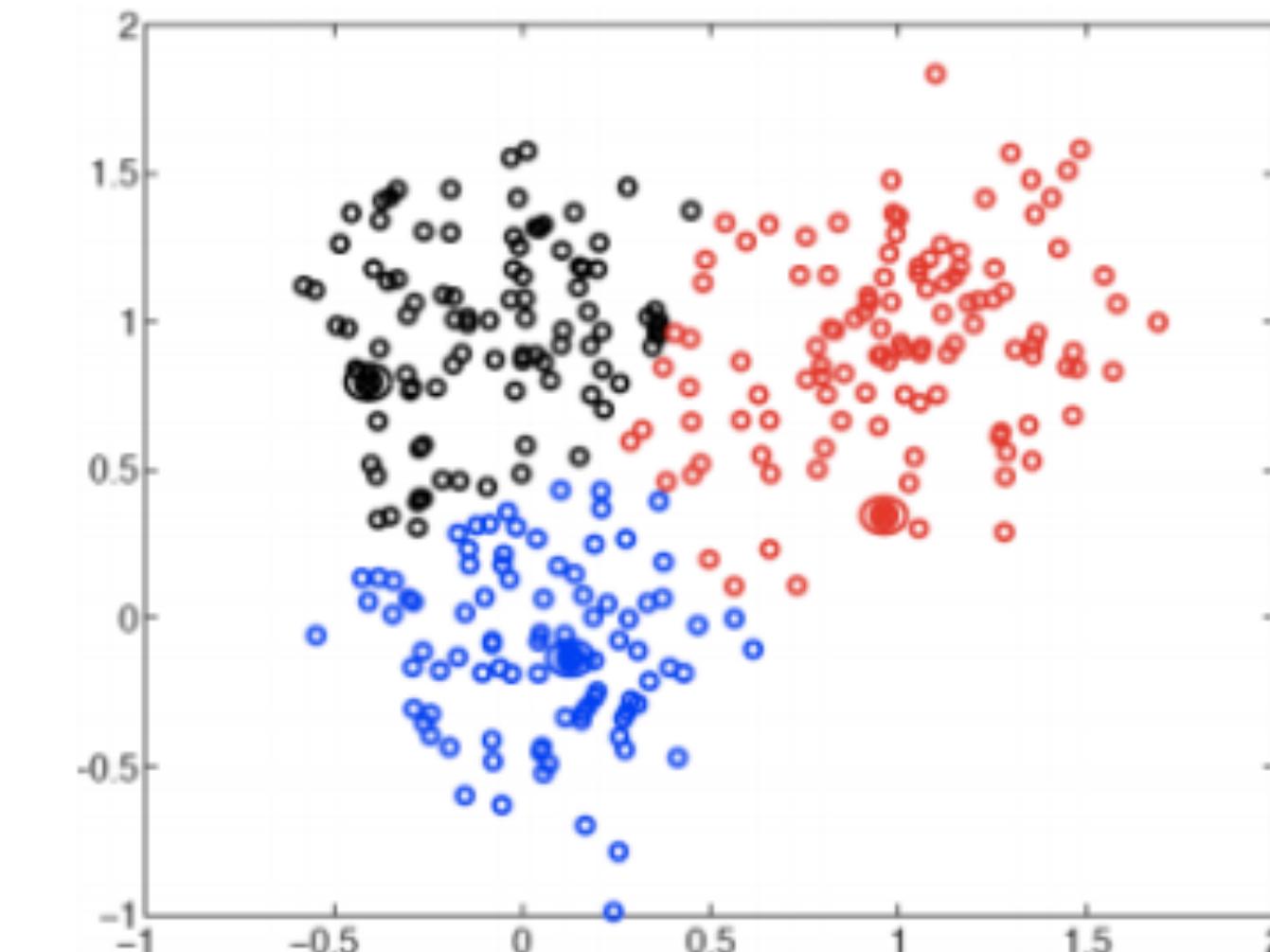
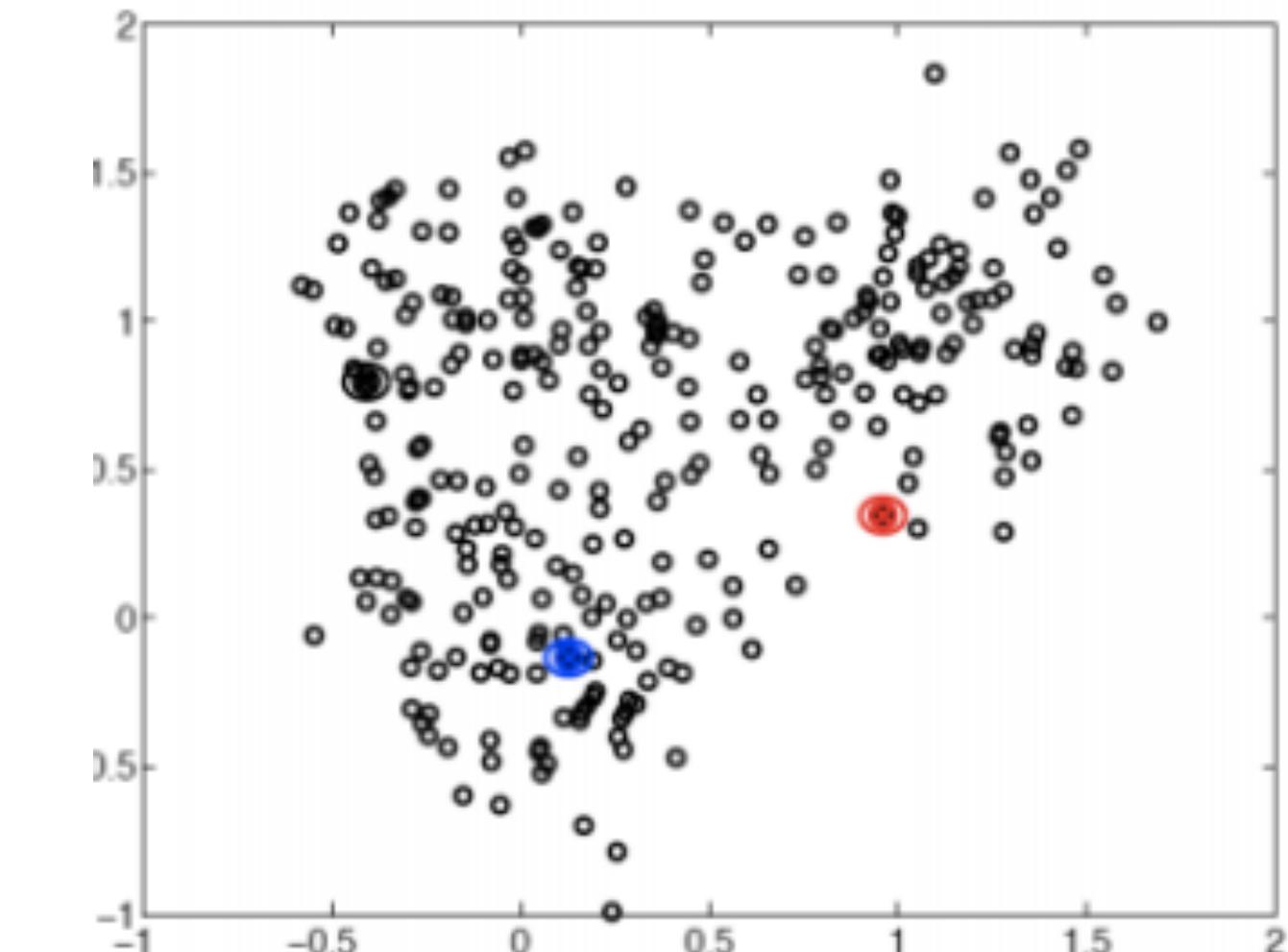
Recap: K-Means Clustering

Iterative procedure

1. Initialization: pick K arbitrary centroids (cluster means)
2. Assign each sample to the closest centroid.
3. Adjust the centroids to be the means of the samples assigned to them.
4. Go to step 2 (until no change)

Algorithm is guaranteed to converge

- Local optimum
- Final result depends on initial estimate



Recap: EM Algorithm

Expectation-Maximization (EM) Algorithm

- **E-step:** softly assign samples to mixture components

$$\gamma_j(\mathbf{x}_n) \leftarrow \frac{\pi_j \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j)}{\sum_{k=1}^K \pi_k \mathcal{N}(\mathbf{x}_n | \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)} \quad \forall j = 1, \dots, K, \quad n = 1, \dots, N$$

- **M-step:** re-estimate the parameters (separately for each mixture component) based on the soft assignments

$$\hat{N}_j \leftarrow \sum_{n=1}^N \gamma_j(\mathbf{x}_n) = \text{soft number of samples labeled } j$$

$$\hat{\pi}_j^{new} \leftarrow \frac{\hat{N}_j}{N}$$

$$\hat{\mu}_j^{new} \leftarrow \frac{1}{\hat{N}_j} \sum_{n=1}^N \gamma_j(\mathbf{x}_n) \mathbf{x}_n$$

$$\hat{\Sigma}_j^{new} \leftarrow \frac{1}{\hat{N}_j} \sum_{n=1}^N \gamma_j(\mathbf{x}_n) (\mathbf{x}_n - \hat{\mu}_j^{new})(\mathbf{x}_n - \hat{\mu}_j^{new})^T$$

Part 1, Video LinearDiscriminantFunction_p1

- Definition
- Extension to multiple classes

Today's Topics

Linear discriminant function

- Definition
- Extension to multiple classes

Least-squares classification

- Derivation
- Shortcomings

Generalized linear models

- Connection to neural networks
- Generalized linear discriminants & gradient descent

Discriminant Functions

Bayesian Decision Theory

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

- Model condition probability densities $p(x|C_k)$ and priors $p(C_k)$
- Compute posteriors $p(C_k|x)$ (using Bayes' rule)
- Minimize probability of misclassification by maximising $p(C|x)$

New Approach

- Directly encode decision boundary
- Without explicit modelling of probability densities
- Minimize misclassification probability directly

Recap: Discriminant Functions

Formulate classification in terms of comparisons

- Discriminant functions

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

- Classify x as class C_k if

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

Examples (Bayes Decision Theory)

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

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

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

Discriminant Functions

Example: 2 classes

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

Decision function (from Bayes Decision Theory)

$$y(x) = p(C_1|x) - p(C_2|x)$$

$$y(x) \models \ln \frac{p(x|C_1)}{p(x|C_2)} + \ln \frac{p(C_1)}{p(C_2)}$$

Learning Discriminant Functions

General classification problem

- Goal: take a new input x and assign it to one of K classes C_k
- Given: training set $X = x_1, \dots, x_N$
with target values $T = t_1, \dots, t_N$
 \Rightarrow Learn a discriminant function $y(x)$ to perform the classification.

2-class problem

- Binary target values: $t_n \in \{0, 1\}$

K -class problem

- 1-of- K coding scheme, e.g. $t_n = (0, 1, 0, 0, 0)^T$

Learning Discriminant Functions

2-class problem

- $y(x) > 0$: Decide for class C_1 , else for class C_2

In the following, we focus on **linear discriminant functions**

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

weight vector “bias” (=threshold)

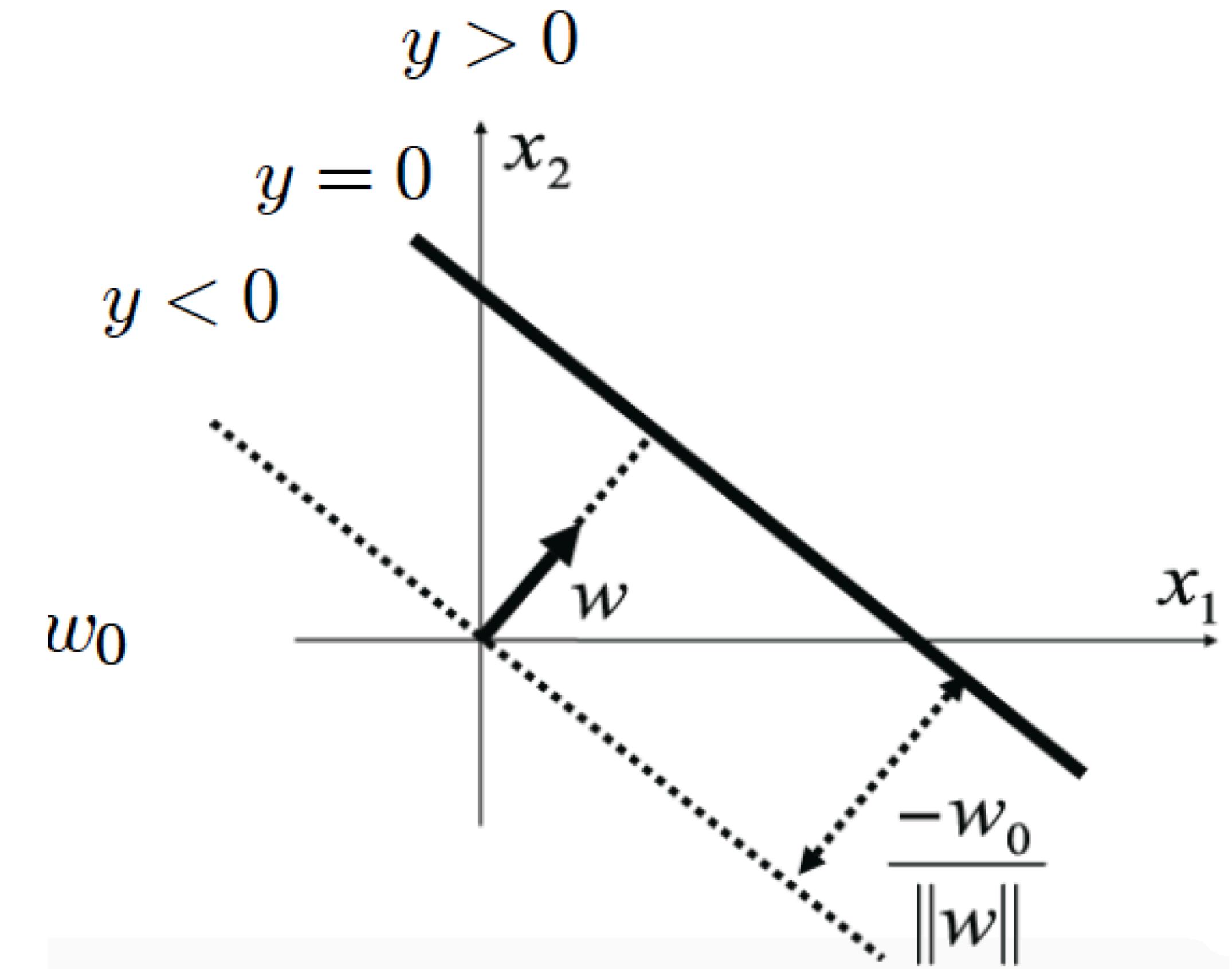
- If a data set can be perfectly classified by a linear discriminant, then we call it **linearly separable**.

Learning Discriminant Functions

Decision boundary $y(x)=0$ defines a hyperplane

- Normal vector: w
- Offset: $\frac{-w_0}{\|w\|}$

$$y(x) = w^T x + w_0$$



Linear Discriminant Functions

Notation

- D : Number of dimensions

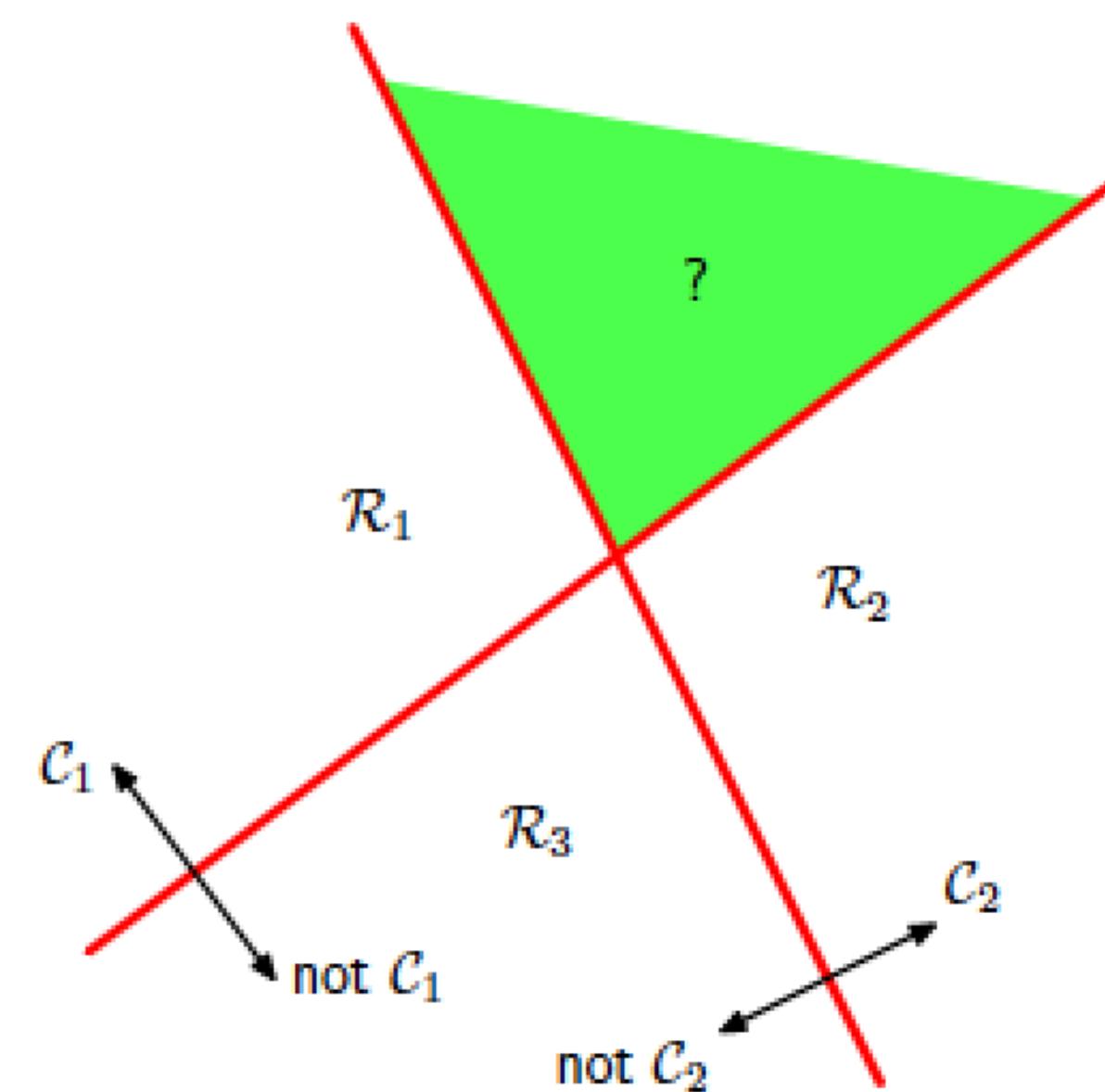
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

$$\begin{aligned} y(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + w_0 \\ &= \sum_{i=1}^D w_i x_i + w_0 \\ &= \sum_{i=0}^D w_i x_i \quad \text{with } x_0 = 1 \text{ constant} \end{aligned}$$

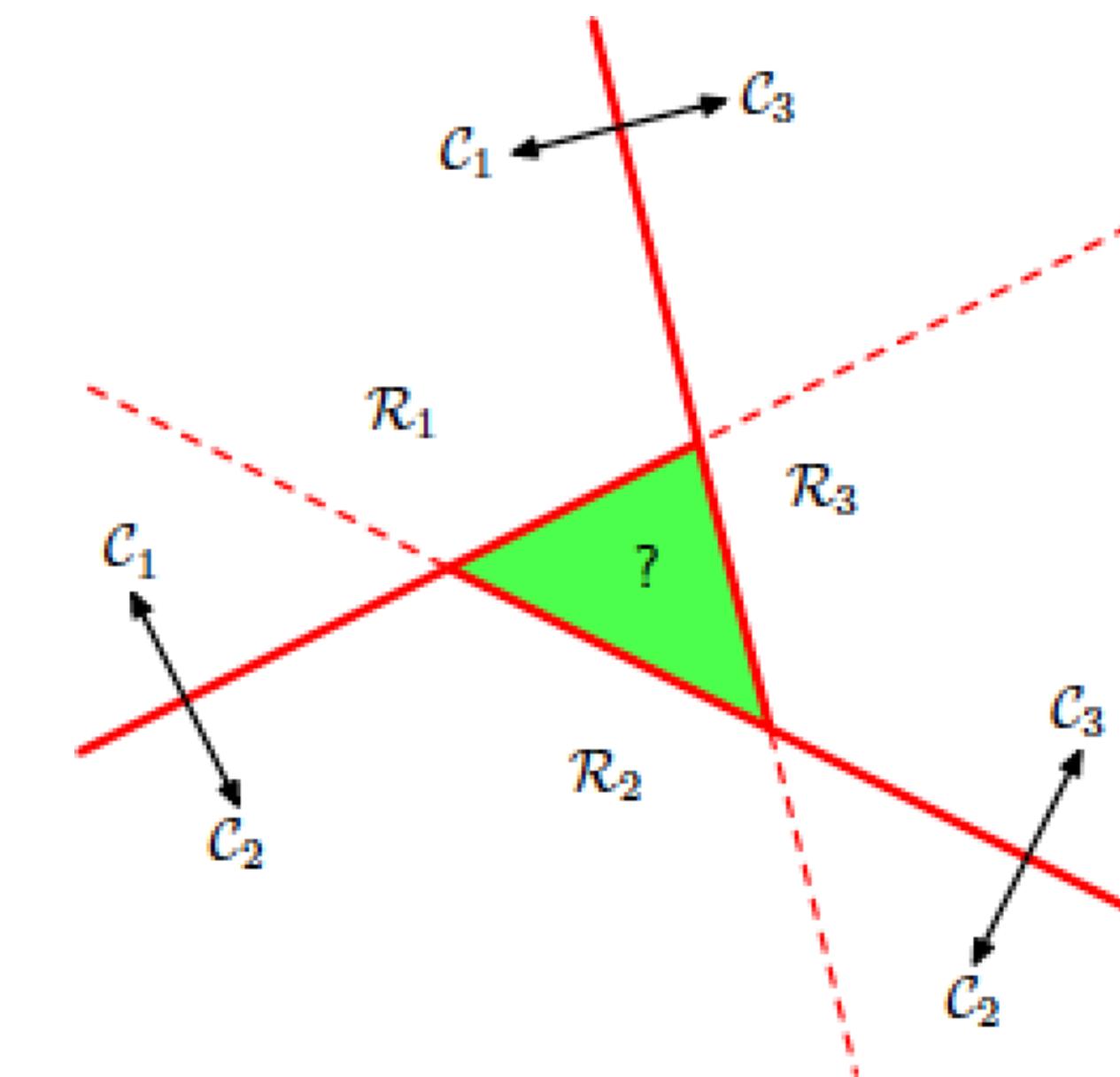
Extension to Multiple Classes

Two simple strategies

- *One-vs-all* classifiers



- *One-vs-one* classifiers



- What difficulties do you see for those strategies?

Extension to Multiple Classes

Problem

- Both strategies in region for which the pure classification result ($y_k > 0$) is ambiguous
- In the *one-vs-all* case, it is still possible to classify those inputs based on the continuous classifier output $y_k > y_j \forall j \neq k$

Solution

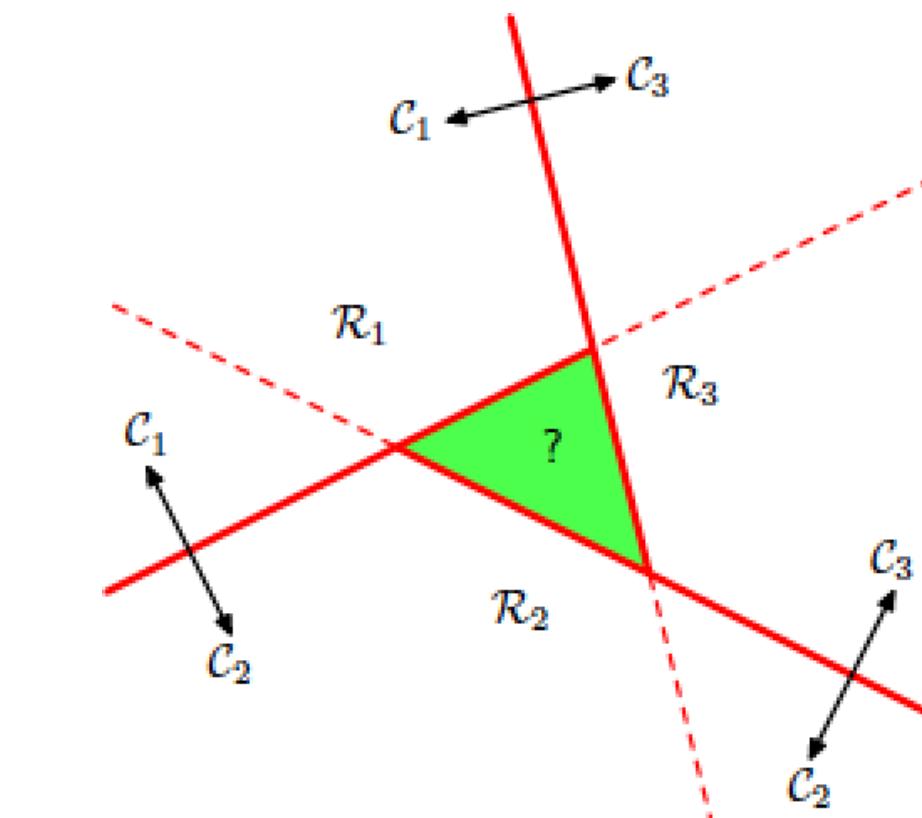
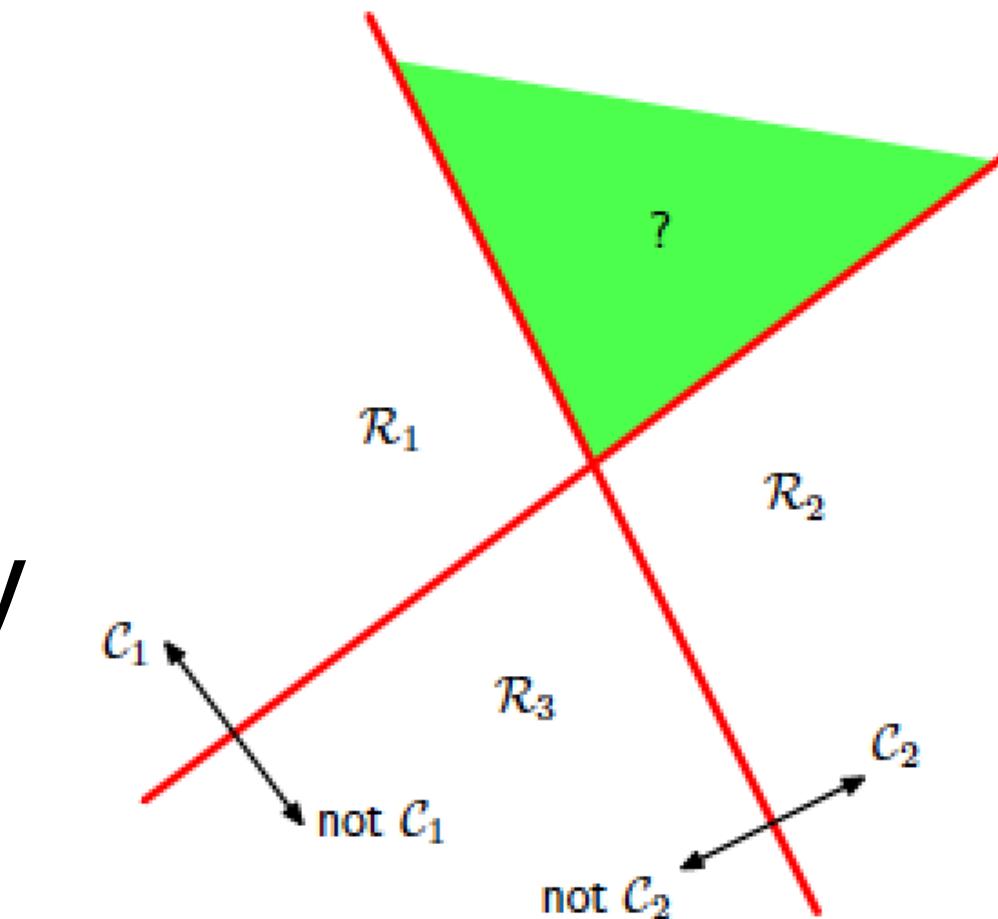
- We can avoid those difficulties by taking K linear functions of the form

$$y_k(x) = w_k^T x + w_{k0}$$

and defining the decision boundaries directly by deciding for C_k iff $y_k > y_j \forall j \neq k$

This corresponds to a 1-of-K coding scheme

$$t_n = (0, 1, 0, 0, 0)^T$$



Extension to Multiple Classes

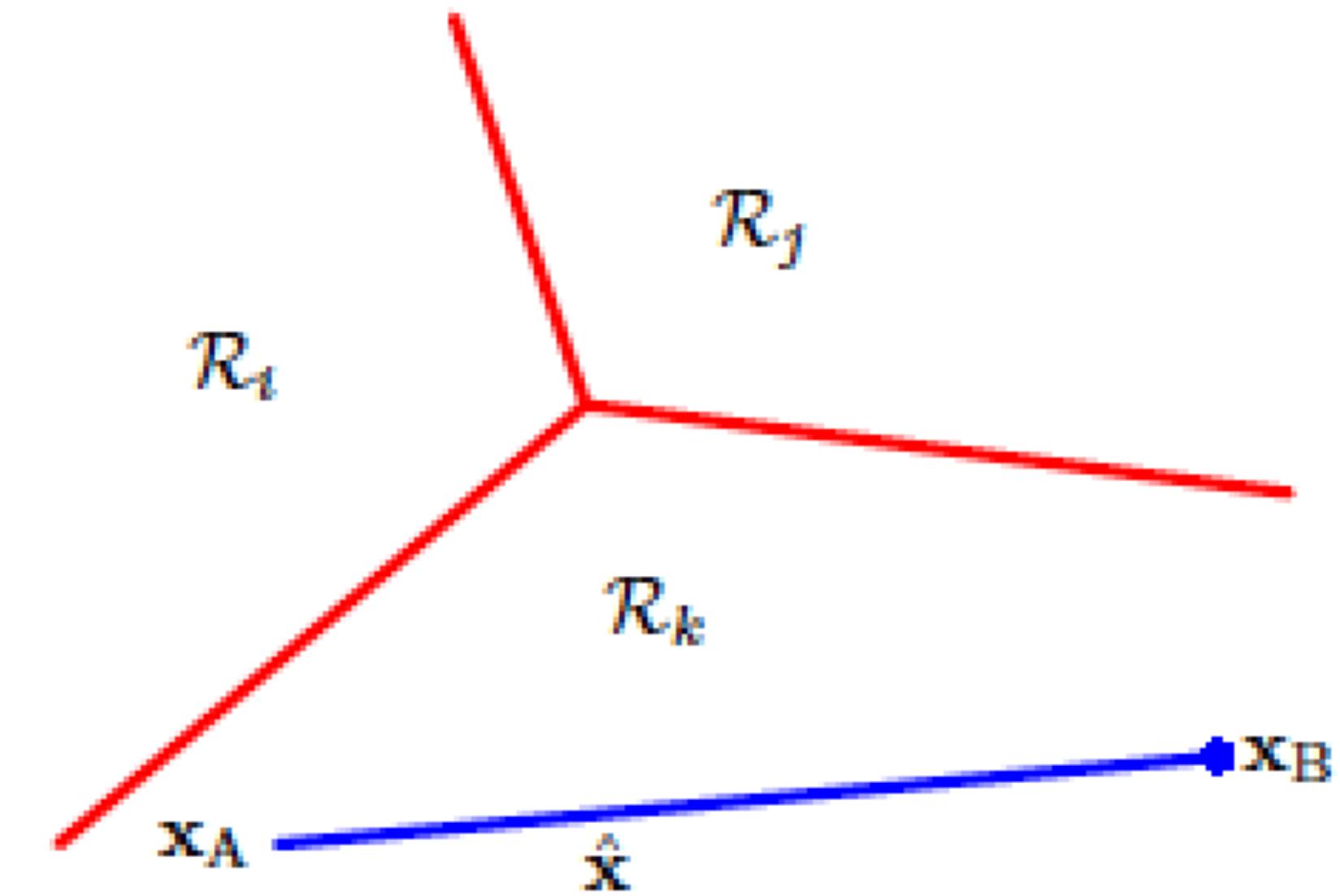
K-class discriminant

- Combination of K linear functions

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

- Resulting decision hyperplanes:

$$(\mathbf{w}_k - \mathbf{w}_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$



- It can be shown that the decision regions of such a discriminant are always singly connected and convex.
- This makes linear discriminant models particularly suitable for problems for which the conditional densities $p(\mathbf{x}|\mathbf{w}_i)$ are unimodal.

Part 2, Video LinearDiscriminantFunction_p2

- Derivation of Least-squares classification
- Shortcomings

Today's Topics

Linear discriminant function

- Definition
- Extension to multiple classes

Least-squares classification

- Derivation
- Shortcomings

Generalized linear models

- Connection to neural networks
- Generalized linear discriminants & gradient descent

General Classification Problem

Classification Problem

- Let's consider K classes described by linear models

$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}, \quad k = 1, \dots, K$$

- We can group those together using vector notation

$$\mathbf{y}(\mathbf{x}) = \tilde{\mathbf{W}}^T \tilde{\mathbf{x}}$$

where

$$\tilde{\mathbf{W}} = [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K] = \begin{bmatrix} w_{10} & \dots & w_{K0} \\ w_{11} & \dots & w_{K1} \\ \vdots & \ddots & \vdots \\ w_{1D} & \dots & w_{KD} \end{bmatrix}$$

- The output will again be in 1-of-K notation.
→ We can directly compare it to the target value $\mathbf{t} = [t_1, \dots, t_k]^T$

General Classification Problem

Classification Problem

- For the entire dataset, we can write

$$\tilde{\mathbf{Y}}(\tilde{\mathbf{X}}) = \tilde{\mathbf{X}}\tilde{\mathbf{W}}$$

and compare to the target matrix \mathbf{T} where

$$\begin{aligned}\tilde{\mathbf{W}} &= [\tilde{\mathbf{w}}_1, \dots, \tilde{\mathbf{w}}_K] \\ \tilde{\mathbf{X}} &= \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad \mathbf{T} = \begin{bmatrix} \mathbf{t}_1^T \\ \vdots \\ \mathbf{t}_N^T \end{bmatrix}\end{aligned}$$

- Result of the comparison: $\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}$

Goal: Choose $\tilde{\mathbf{W}}$ such that this is minimal!

Least-Squares Classification

Simplest Approach

- Directly try to minimise the **sum-of-squares error**
- We could write this as

$$\begin{aligned} E(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn})^2 \\ &= \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (\mathbf{w}_k^T \mathbf{x}_n - t_{kn})^2 \end{aligned}$$

- But let's stick with the matrix notation for now ...
- (The result will be simpler to express and we'll learn some nice matrix algebra rules along the way ...)

Least-Squares Classification

Multi-class case

- Let's formulate the sum-of-squares error in matrix notation

$$E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

using:
 $\sum_{i,j} a_{ij}^2 = \text{Tr}\{\mathbf{A}^T \mathbf{A}\}$

- Taking the derivative yields

$$\frac{\partial}{\partial \tilde{\mathbf{W}}} E_D(\tilde{\mathbf{W}}) = \frac{1}{2} \frac{\partial}{\partial \tilde{\mathbf{W}}} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

chain rule:

$$\frac{\partial \mathbf{Z}}{\partial \mathbf{X}} = \frac{\partial \mathbf{Z}}{\partial \mathbf{Y}} \frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$$

$$= \frac{1}{2} \frac{\partial}{\partial (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})} \text{Tr} \left\{ (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T}) \right\}$$

$$\cdot \frac{\partial}{\partial \tilde{\mathbf{W}}} (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})$$

$$= \tilde{\mathbf{X}}^T (\tilde{\mathbf{X}}\tilde{\mathbf{W}} - \mathbf{T})$$

using:

$$\frac{\partial}{\partial \mathbf{A}} \text{Tr} \{ \mathbf{A} \} = \mathbf{I}$$

Least-Squares Classification

Minimizing the sum-of-squares error

$$\frac{\partial}{\partial \widetilde{\mathbf{W}}} E_D(\widetilde{\mathbf{W}}) = \widetilde{\mathbf{X}}^T (\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} - \mathbf{T}) \stackrel{!}{=} 0$$

$$\widetilde{\mathbf{X}} \widetilde{\mathbf{W}} = \mathbf{T}$$

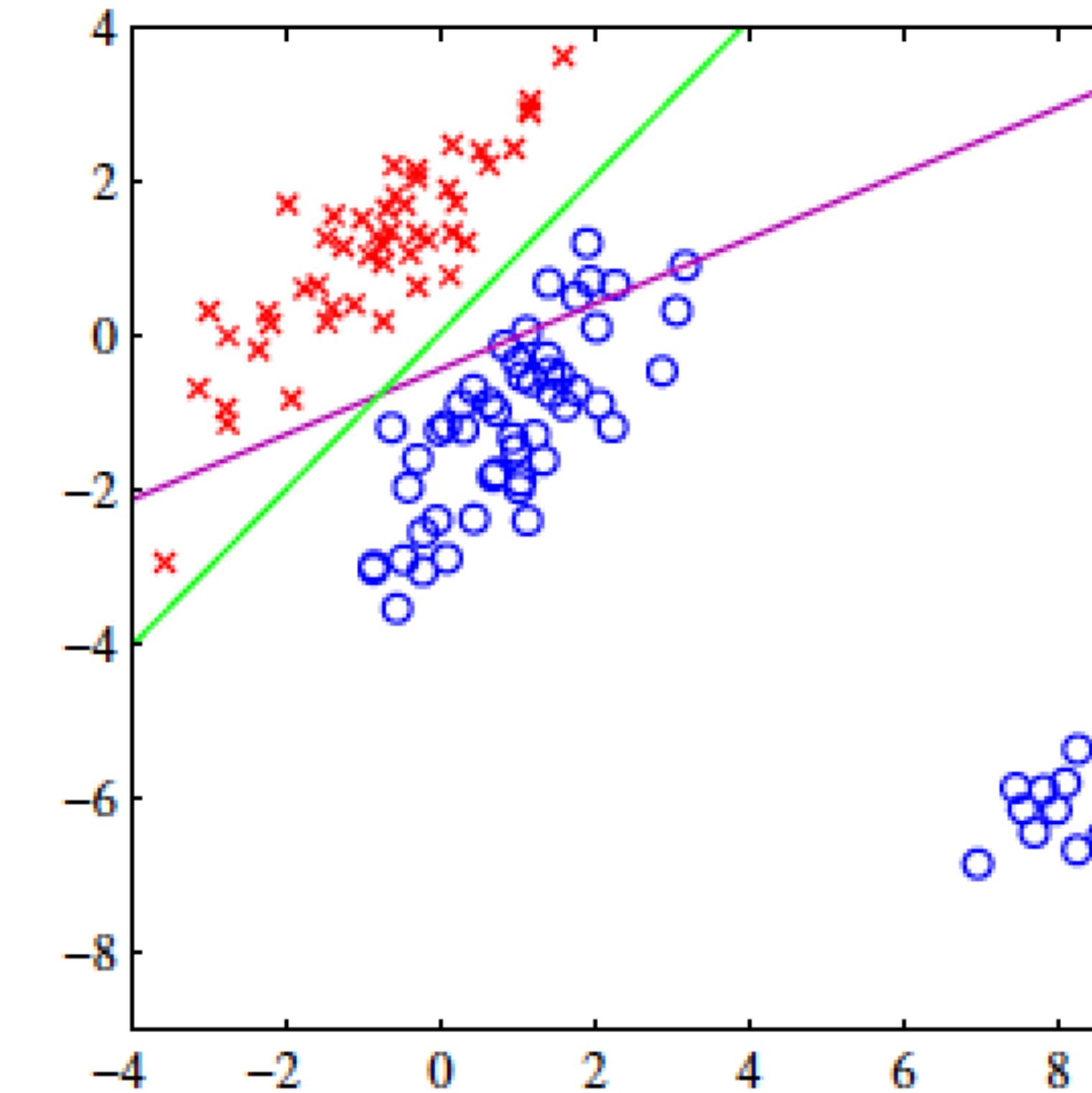
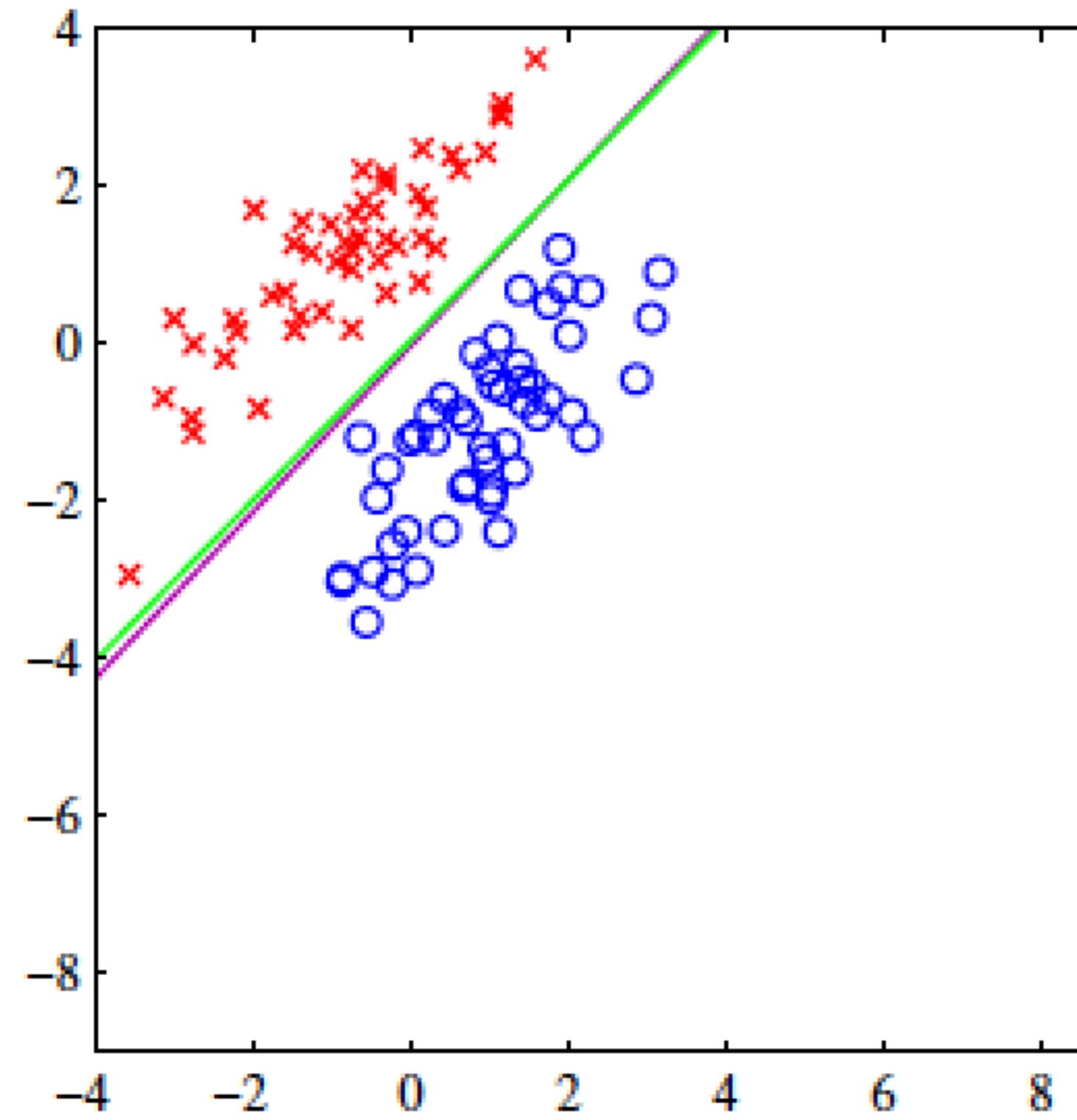
$$\begin{aligned}\widetilde{\mathbf{W}} &= (\widetilde{\mathbf{X}}^T \widetilde{\mathbf{X}})^{-1} \widetilde{\mathbf{X}}^T \mathbf{T} \\ &= \widetilde{\mathbf{X}}^\dagger \mathbf{T} \quad \text{“pseudo-inverse”}\end{aligned}$$

- We then obtain the discriminant function as

$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}} = \mathbf{T}^T (\widetilde{\mathbf{X}}^\dagger)^T \widetilde{\mathbf{x}}$$

- Exact, closed-form solution for the discriminant function parameters

Problems with Least Squares



Least-squares is very sensitive to outliers!

- The error function penalises predictions that are “too correct”

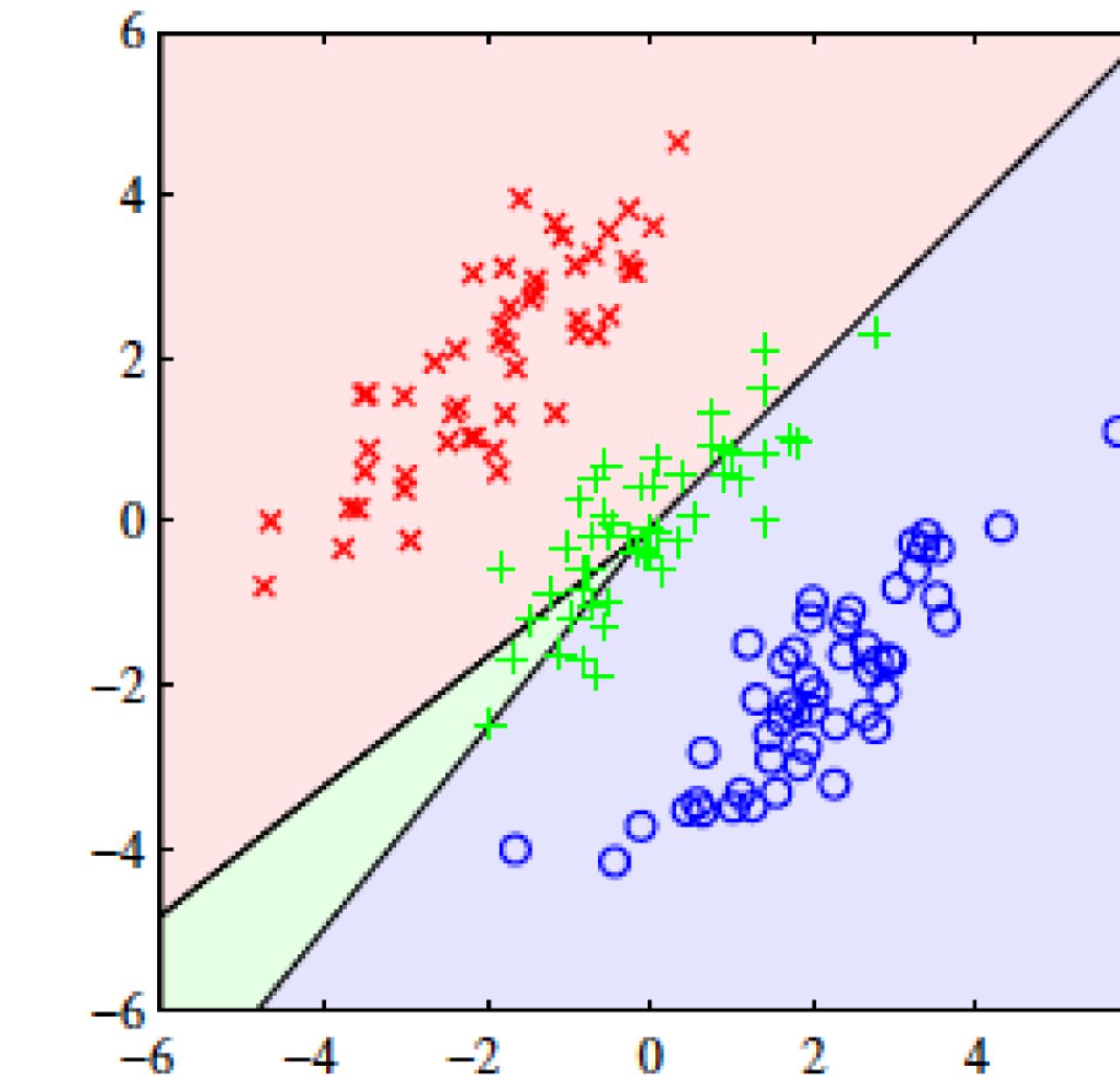
Problems with Least Squares

Another example:

- 3 classes (red, green, blue)
- Linearly separable problem
- Least-squares solution:
Most green points are misclassified!

Deeper reason for the failure

- Least-squares corresponds to Maximum Likelihood under the assumption of a Gaussian conditional distribution.
 - However, our binary target vectors have a distribution that is clearly non-Gaussian!
- Least-squares is the wrong probabilistic tool in this case!



Part 3, Video LinearDiscriminantFunction_p3

- Connection to neural networks
- Generalized linear discriminants & gradient descent

Today's Topics

Linear discriminant function

- Definition
- Extension to multiple classes

Least-squares classification

- Derivation
- Shortcomings

Generalized linear models

- Connection to neural networks
- Generalized linear discriminants & gradient descent

Generalized Linear Models

Linear model

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

Generalized linear model

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

- $g(\cdot)$ is called an **activation function** and maybe nonlinear.
- The decision surface correspond to

$$y(\mathbf{x}) = \text{const.} \Leftrightarrow \mathbf{w}^T \mathbf{x} + w_0 = \text{const.}$$

- If g is monotonous (which is typically the case), the resulting decision boundaries are still linear functions of \mathbf{x} .

Generalized Linear Models

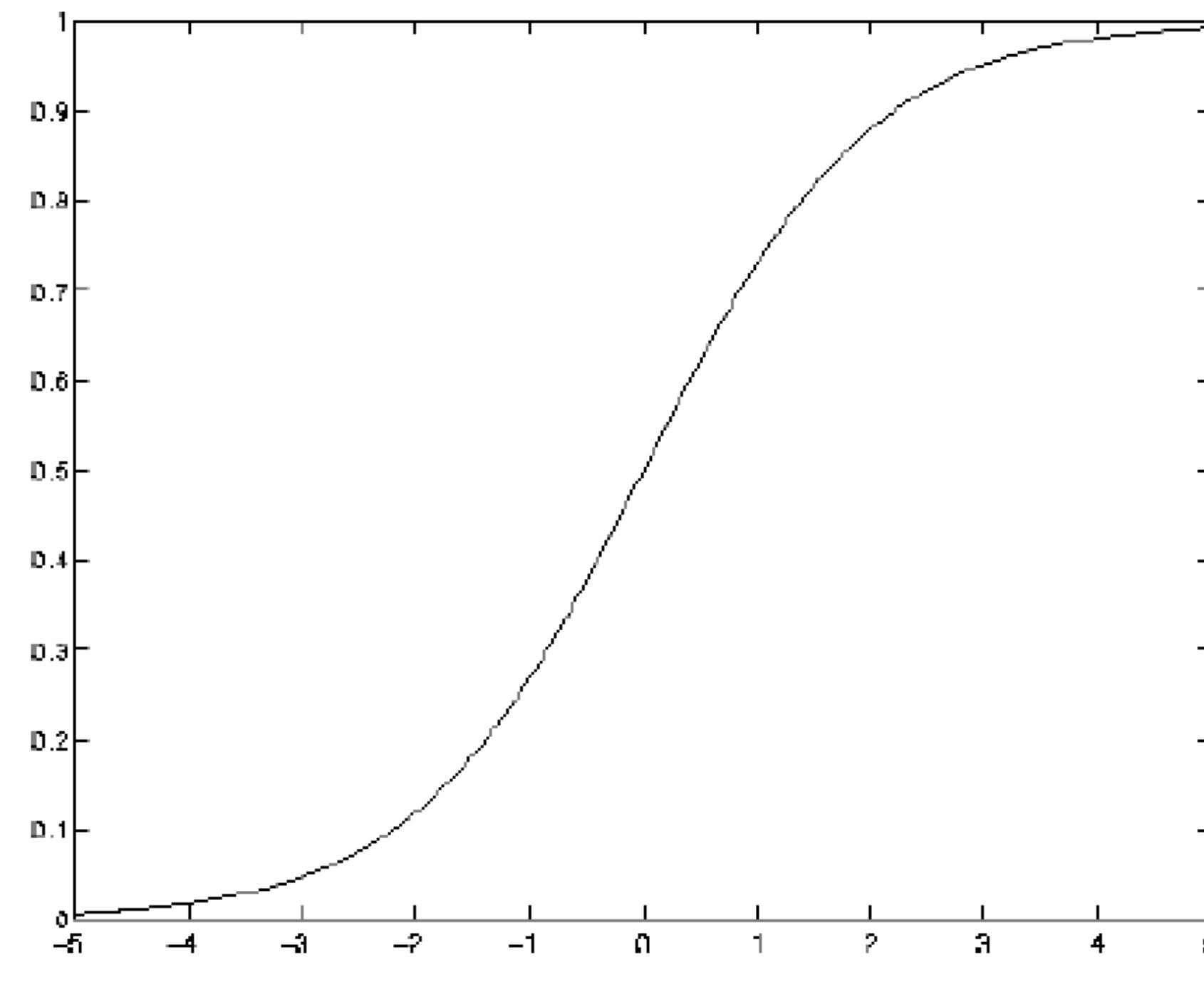
Consider 2 classes:

$$\begin{aligned} p(\mathcal{C}_1|\mathbf{x}) &= \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)} \\ &= \frac{1}{1 + \frac{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}} \\ &= \frac{1}{1 + \exp(-a)} \equiv g(a) \end{aligned}$$

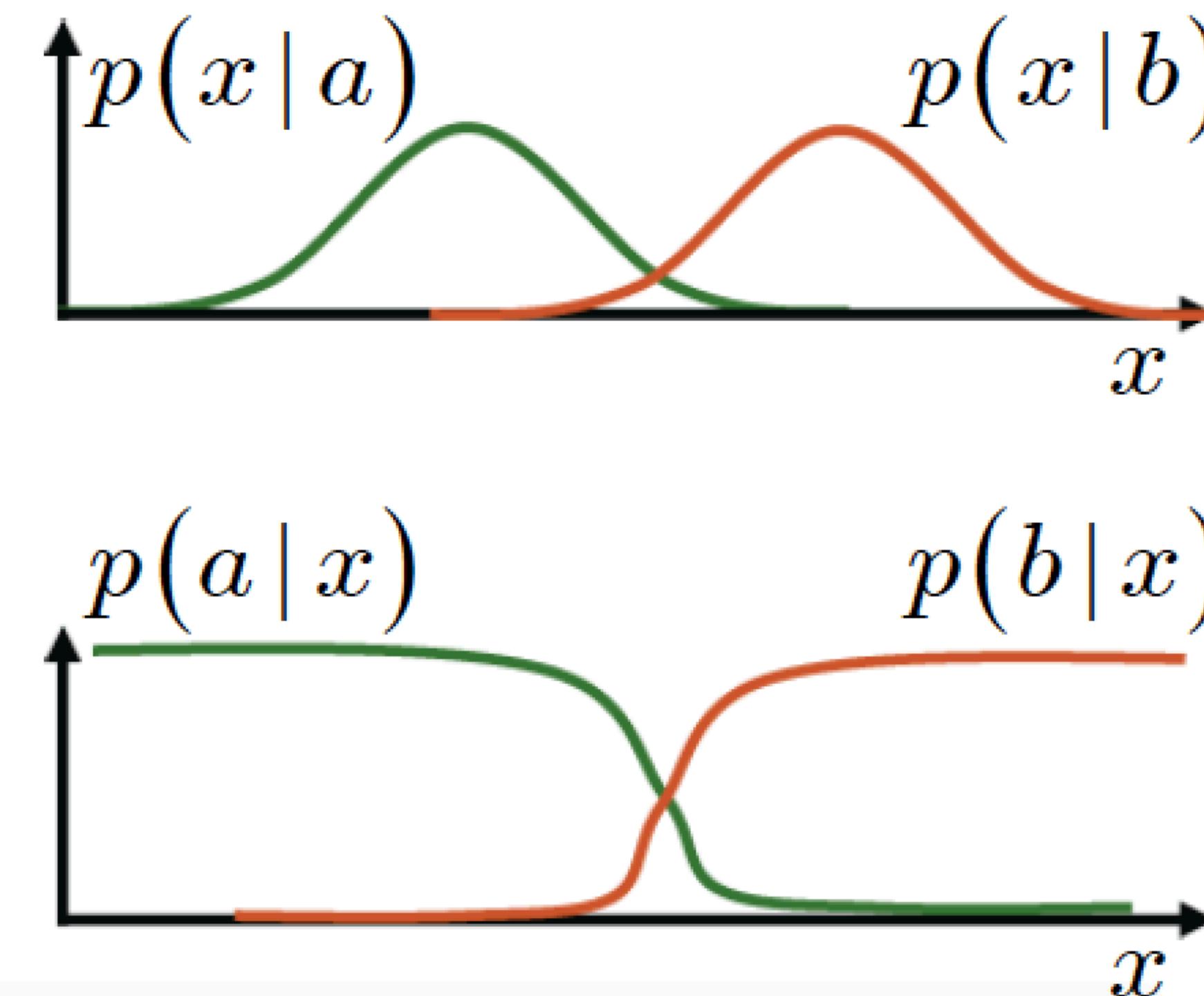
$$\text{with } a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$

Logistic Sigmoid Activation Function

$$g(a) \equiv \frac{1}{1 + \exp(-a)}$$



Example: Normal distributions
with identical covariance



Normalized Exponential

General case of K>2 classes:

$$\begin{aligned} p(\mathcal{C}_k | \mathbf{x}) &= \frac{p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j) p(\mathcal{C}_j)} \\ &= \frac{\exp(a_k)}{\sum_j \exp(a_j)} \end{aligned}$$

$$\text{with } a_k = \ln p(\mathbf{x} | \mathcal{C}_k) p(\mathcal{C}_k)$$

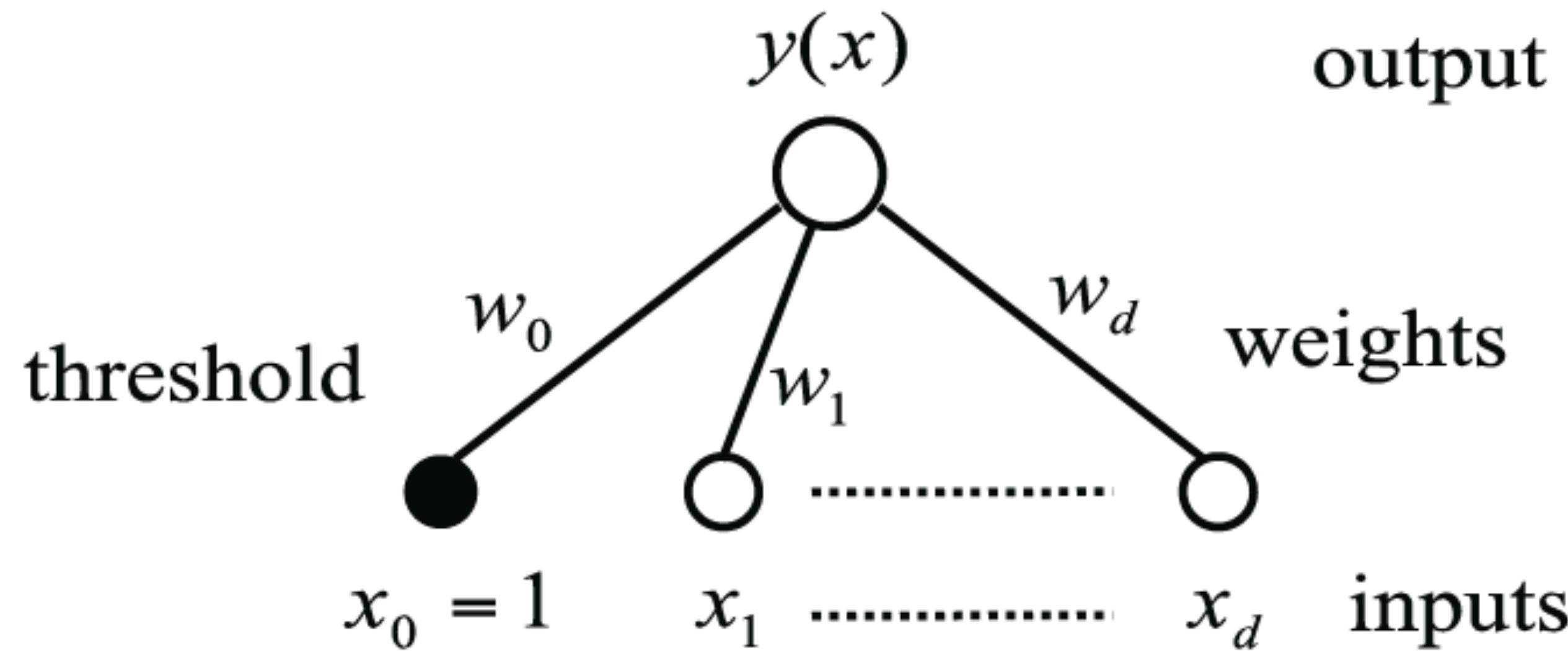
- This is known as the **normalized exponential** or **softmax** function
- Can be regarded as a multi class generalised of the logistic sigmoid

Relationship to Neural Networks

2-Class case

$$y(x) = g\left(\sum_{i=0}^D w_i x_i\right) \text{ with } x_0 = 1 \text{ constant}$$

Neural network (“single-layer perceptron”)

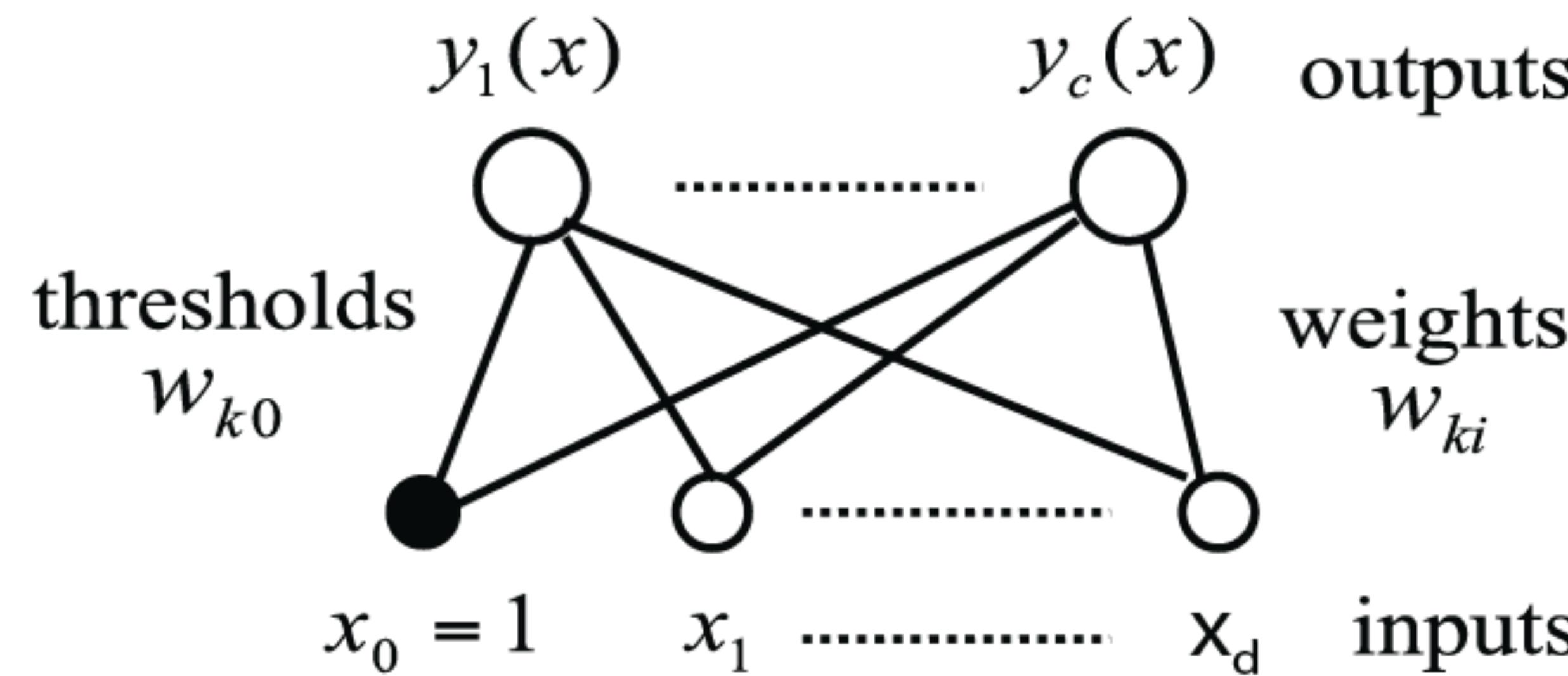


Relationship to Neural Networks

Multi-class case

$$y_k(x) = g\left(\sum_{i=0}^D w_{ki}x_i\right) \text{ with } x_0 = 1 \text{ constant}$$

Multi-class perceptron

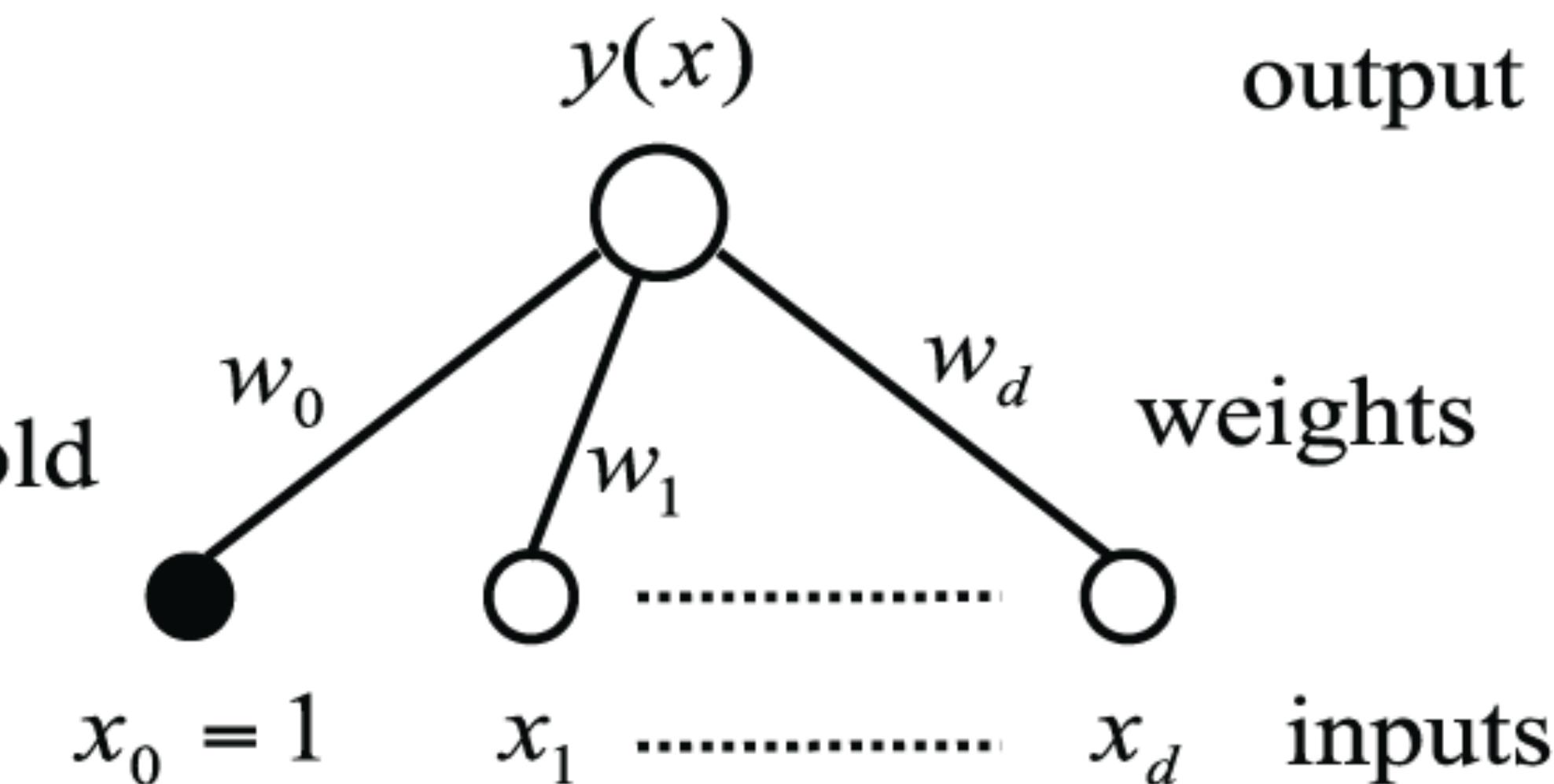


Logistic Discrimination

If we use the logistic sigmoid activation function...

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

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



... then we can interpret the $y(x)$ as posterior probabilities!

Other Motivation for Nonlinearity

Recall least-squares classification

- One of the problems was that data points that are “too correct” have a strong influence on the decision surface under a squared-error criterion.

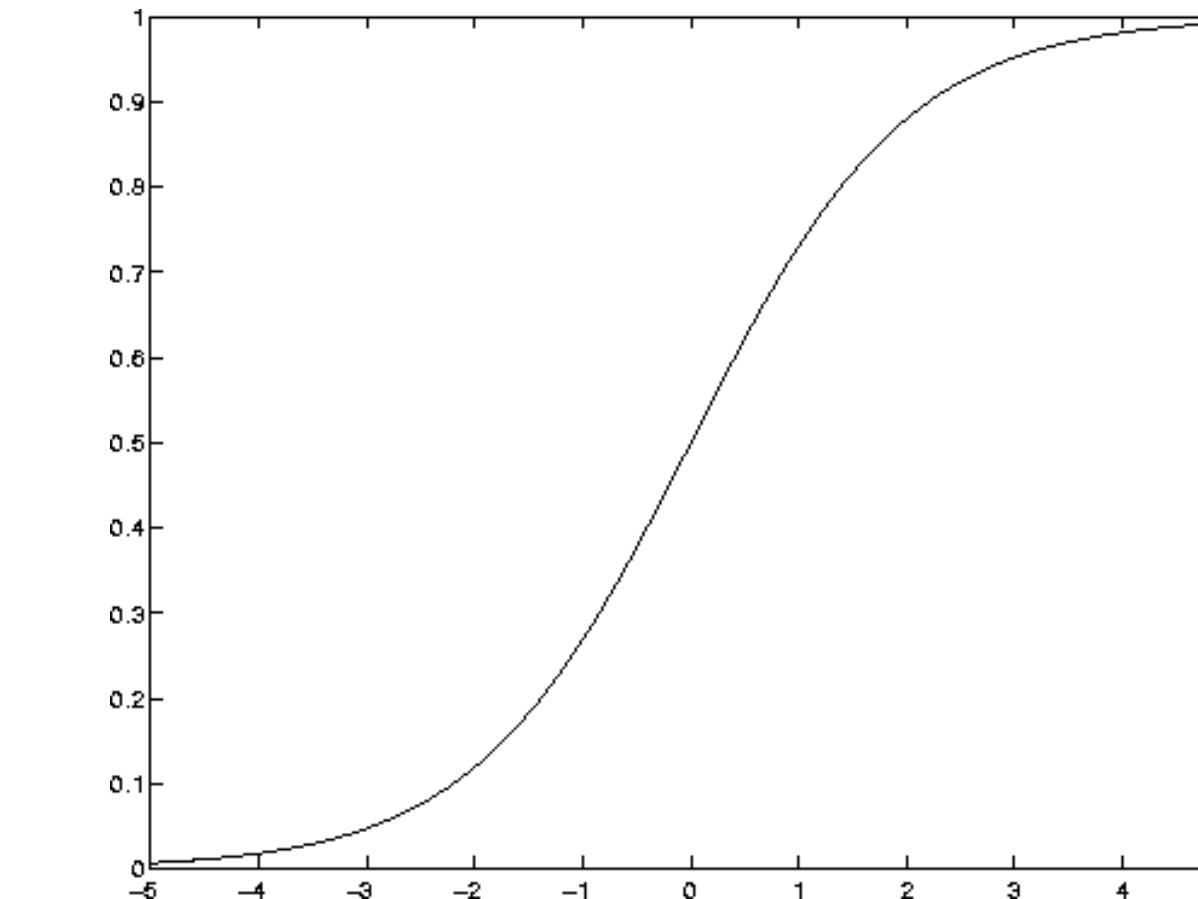
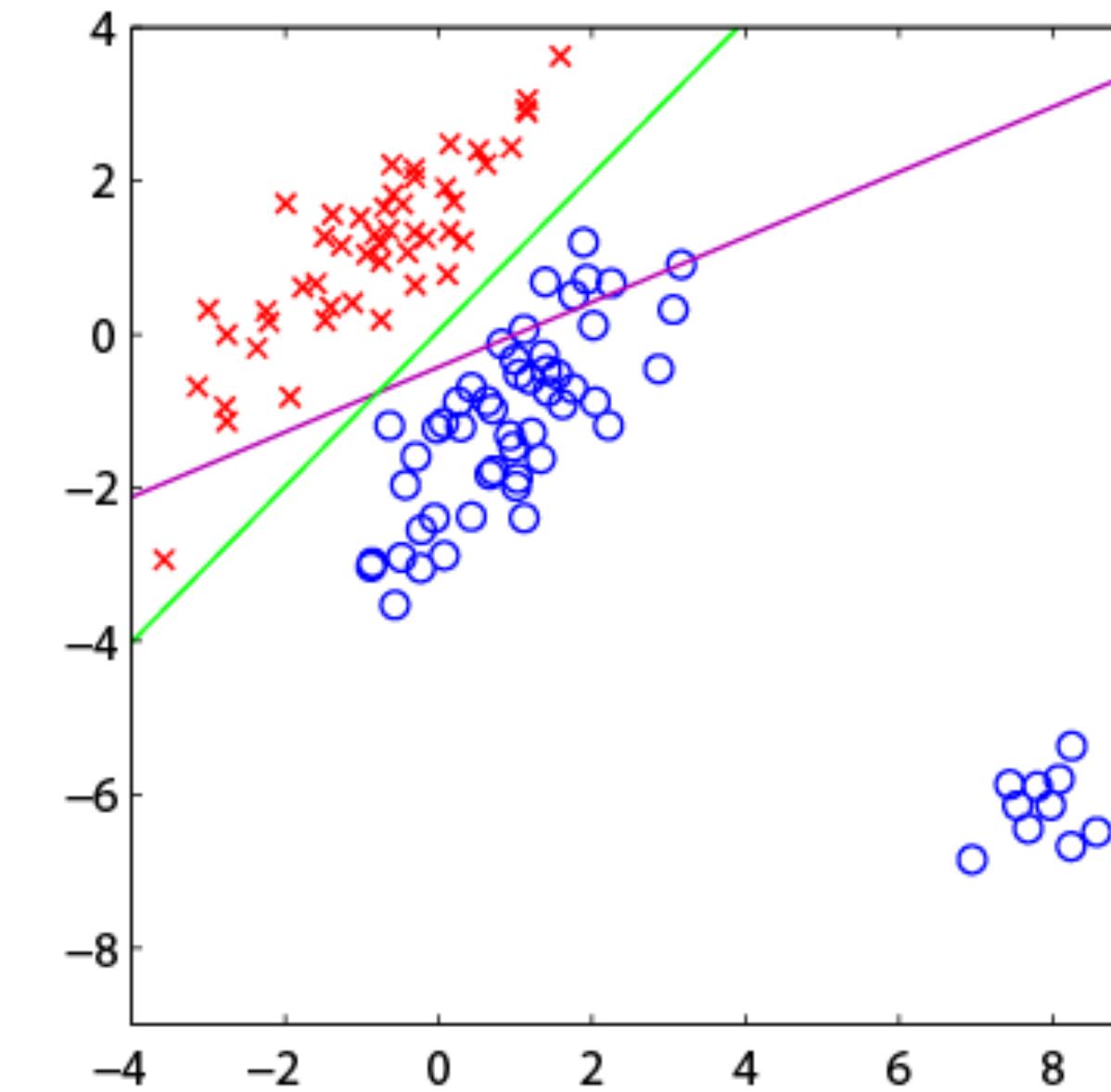
$$E(\mathbf{w}) = \sum_{n=1}^N (y(\mathbf{x}_n; \mathbf{w}) - t_n)^2$$

- Reason: the output of $y(\mathbf{x}_n; \mathbf{w})$ can grow arbitrarily large for some \mathbf{x}_n :

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

- By choosing a suitable nonlinearity (e.g. a sigmoid), we can limit those influences

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



Discussion: Generalized Linear Models

Advantages

- The nonlinearity gives us more flexibility.
- Can be used to limit the effect of outliers
- Choice of a sigmoid leads to a nice probabilistic interpretation

Disadvantages

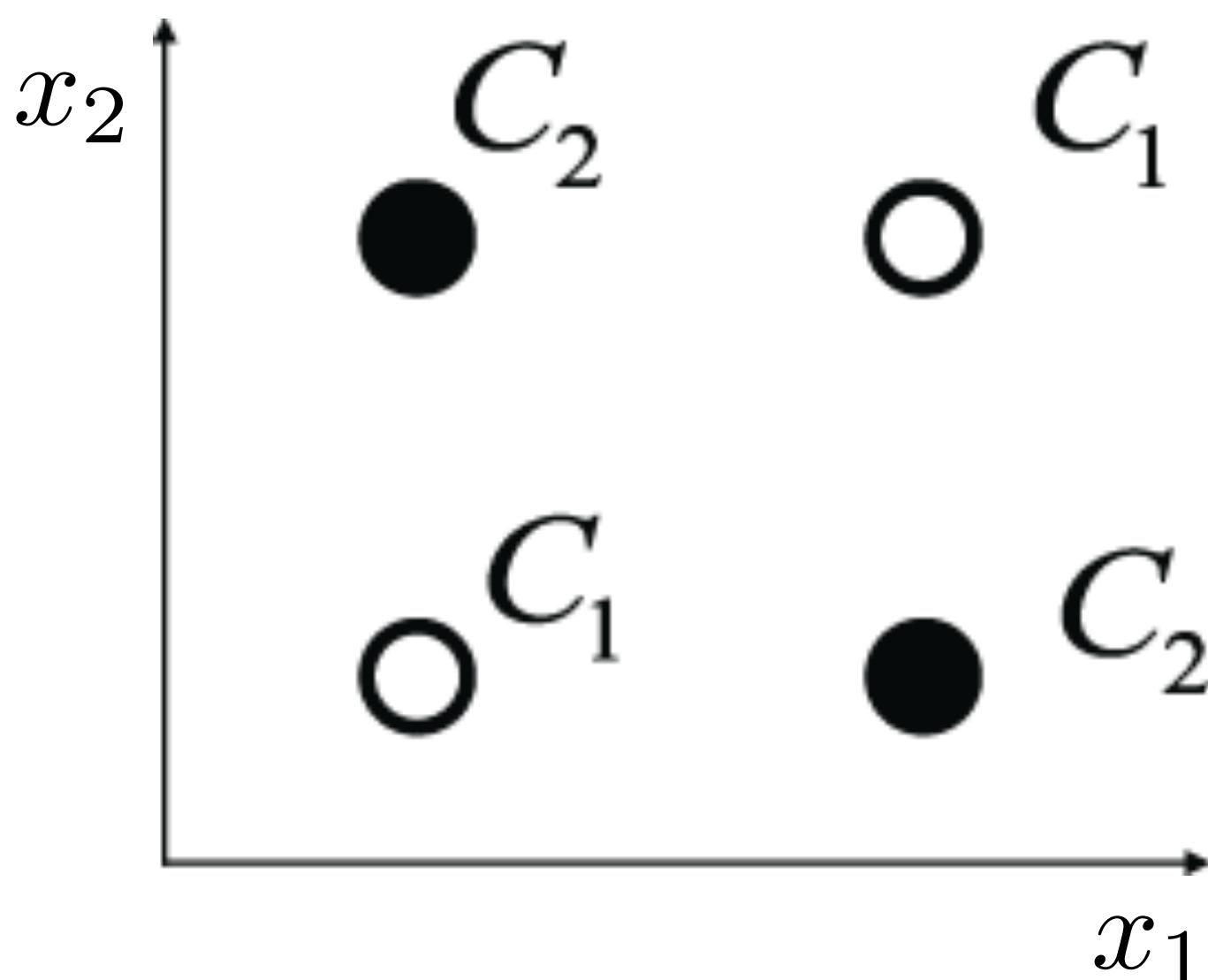
- Least-squares minimisation in general no longer leads to a close-form analytical solution.
- Need to apply iterative methods.
→ Gradient descent.

Linear Separability

Up to now: restrictive assumption

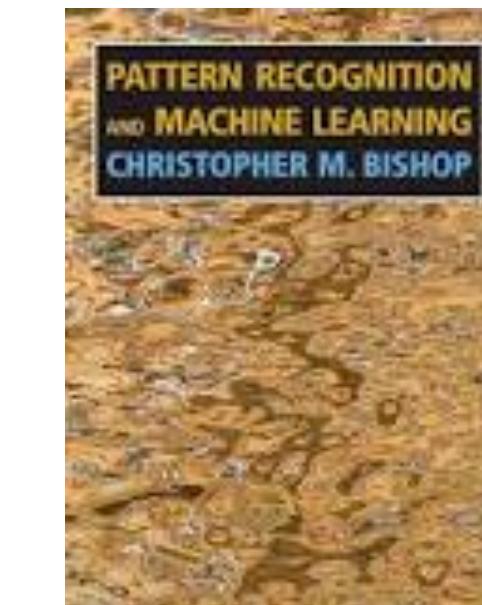
- Only consider linear decision boundaries

Classical counterexample: XOR



Readings

Bishop's book



More information on Linear Discriminant Functions can be found in Chapter 4 of Bishop's book (in particular Chapter 4.1).

You should be able to answer ...

- What is 1 of K coding scheme?
- Formula for linear discriminant function, explain its entries.
- What does linearly separable mean?
- Extension of linear discriminant to multiple classes - two strategies.
What are the difficulties with these strategies?
- What is a sum-of-squares error?
- Problems with least squares?
- Difference between linear model and generalised linear model.
- Sigmoid activation function: formula and how to implement it.
- What is softmax?
- What are advantages and disadvantages of generalised linear model?