

RWTH AACHEN  
UNIVERSITY

# Machine Learning – Lecture 5

## Linear Discriminant Functions

23.10.2019

Bastian Leibe  
RWTH Aachen  
<http://www.vision.rwth-aachen.de>  
leibe@vision.rwth-aachen.de

Machine Learning Winter '19

RWTH AACHEN  
UNIVERSITY

## Course Outline

- Fundamentals
  - Bayes Decision Theory
  - Probability Density Estimation
- Classification Approaches
  - Linear Discriminants
  - Support Vector Machines
  - Ensemble Methods & Boosting
  - Randomized Trees, Forests & Ferns
- Deep Learning
  - Foundations
  - Convolutional Neural Networks
  - Recurrent Neural Networks

B. Leibe

RWTH AACHEN  
UNIVERSITY

## Recap: Mixture of Gaussians (MoG)

- “Generative model”

$p(j) = \pi_j$  “Weight” of mixture component

$p(x|\theta_j)$  Mixture component

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

Slide credit: Bernt Schiele

RWTH AACHEN  
UNIVERSITY

## Recap: Estimating MoGs – Iterative Strategy

- Assuming we knew the values of the hidden variable...

assumed known	→ 1 111	22 2 2	j
$h(j=1 x_n)$	1 111	00 0 0	
$h(j=2 x_n)$	0 000	11 1 1	

$$\mu_1 = \frac{\sum_{n=1}^N h(j=1|x_n)x_n}{\sum_{i=1}^N h(j=1|x_n)} \quad \mu_2 = \frac{\sum_{n=1}^N h(j=2|x_n)x_n}{\sum_{i=1}^N h(j=2|x_n)}$$

Slide credit: Bernt Schiele

RWTH AACHEN  
UNIVERSITY

## Recap: Estimating MoGs – Iterative Strategy

- Assuming we knew the mixture components...

1 111	22 2 2	j
-------	--------	---

- Bayes decision rule: Decide  $j = 1$  if
 
$$p(j=1|x_n) > p(j=2|x_n)$$

Slide credit: Bernt Schiele

RWTH AACHEN  
UNIVERSITY

## Recap: K-Means Clustering

- Iterative procedure
  - Initialization: pick  $K$  arbitrary centroids (cluster means)
  - Assign each sample to the closest centroid.
  - Adjust the centroids to be the means of the samples assigned to them.
  - Go to step 2 (until no change)
- Algorithm is guaranteed to converge after finite #iterations.
  - Local optimum
  - Final result depends on initialization.

Slide credit: Bernt Schiele

RWTH AACHEN UNIVERSITY

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

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

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

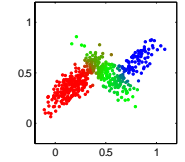
$$\hat{\boldsymbol{\Sigma}}_j^{\text{new}} \leftarrow \frac{1}{\hat{N}_j} \sum_{n=1}^N \gamma_j(\mathbf{x}_n) (\mathbf{x}_n - \hat{\boldsymbol{\mu}}_j^{\text{new}})(\mathbf{x}_n - \hat{\boldsymbol{\mu}}_j^{\text{new}})^T$$

Machine Learning Winter '19 | Slide adapted from Bernt Schiele | B. Leibe | 7

RWTH AACHEN UNIVERSITY

## Applications

- Mixture models are used in many practical applications.
  - Wherever distributions with complex or unknown shapes need to be represented...
- Popular application in Computer Vision
  - Model distributions of pixel colors.
  - Each pixel is one data point in, e.g., RGB space.
    - Learn a MoG to represent the class-conditional densities.
    - Use the learned models to classify other pixels.

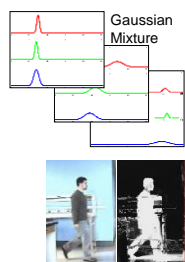


Machine Learning Winter '19 | B. Leibe | Image source: C.M. Bishop, 2006 | 8

RWTH AACHEN UNIVERSITY

## Application: Background Model for Tracking


- Train background MoG for each pixel
  - Model "common" appearance variation for each background pixel.
  - Initialization with an empty scene.
  - Update the mixtures over time
    - Adapt to lighting changes, etc.
- Used in many vision-based tracking applications
  - Anything that cannot be explained by the background model is labeled as foreground (=object).
  - Easy segmentation if camera is fixed.



C. Stauffer, E. Grimson, *Learning Patterns of Activity Using Real-Time Tracking*, IEEE Trans. PAMI, 22(8):747-757, 2000. | B. Leibe | Image Source: Daniel Roth, Tobias Jaeger | 9

RWTH AACHEN UNIVERSITY

## Application: Image Segmentation



- User assisted image segmentation
  - User marks two regions for foreground and background.
  - Learn a MoG model for the color values in each region.
  - Use those models to classify all other pixels.
    - Simple segmentation procedure (building block for more complex applications)

Machine Learning Winter '19 | B. Leibe | 10

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Linear discriminant functions
  - Definition
  - Extension to multiple classes
- Least-squares classification
  - Derivation
  - Shortcomings
- Generalized linear models
  - Connection to neural networks
  - Generalized linear discriminants & gradient descent

Machine Learning Winter '19 | B. Leibe | 11

RWTH AACHEN UNIVERSITY

## Discriminant Functions

- Bayesian Decision Theory
 
$$p(C_k | x) = \frac{p(x | C_k) p(C_k)}{p(x)}$$
  - Model conditional 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 maximizing  $p(C | x)$
- New approach
  - Directly encode decision boundary
  - Without explicit modeling of probability densities
  - Minimize misclassification probability directly.

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 12

RWTH AACHEN  
UNIVERSITY

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

Machine Learning Winter '19 13  
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Discriminant Functions

- Example: 2 classes
  - $$y_1(x) > y_2(x)$$
  - $$\Leftrightarrow y_1(x) - y_2(x) > 0$$
  - $$\Leftrightarrow \mathbf{y}(x) > 0$$
- Decision functions (from Bayes Decision Theory)
  - $$y(x) = p(C_1|x) - p(C_2|x)$$
  - $$y(x) = \ln \frac{p(x|C_1)}{p(x|C_2)} + \ln \frac{p(C_1)}{p(C_2)}$$

Machine Learning Winter '19 14  
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Learning Discriminant Functions

- General classification problem
  - Goal: take a new input  $\mathbf{x}$  and assign it to one of  $K$  classes  $C_k$ .
  - Given: training set  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$  with target values  $\mathbf{T} = \{t_1, \dots, t_N\}$ .
  - $\Rightarrow$  Learn a discriminant function  $y(\mathbf{x})$  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.  $t_n = (0, 1, 0, 0, 0)^T$

Machine Learning Winter '19 15  
B. Leibe

RWTH AACHEN  
UNIVERSITY

## Linear 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**

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$
  - weight vector  $\mathbf{w}$
  - "bias"  $w_0$  (= threshold)
- If a data set can be perfectly classified by a linear discriminant, then we call it **linearly separable**.

Machine Learning Winter '19 16  
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Linear Discriminant Functions

- Decision boundary  $y(\mathbf{x}) = 0$  defines a hyperplane
  - Normal vector:  $\mathbf{w}$
  - Offset:  $\frac{-w_0}{\|\mathbf{w}\|}$

Machine Learning Winter '19 17  
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN  
UNIVERSITY

## Linear Discriminant Functions

- Notation
  - $D$  : Number of dimensions
$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_D \end{bmatrix} \quad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_D \end{bmatrix}$$

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

$$= \sum_{i=1}^D w_i x_i + w_0$$

$$= \sum_{i=0}^D w_i x_i \quad \text{with } x_0 = 1 \text{ constant}$$

Machine Learning Winter '19 18  
Slide credit: Bernt Schiele B. Leibe

Machine Learning Winter '19

## Extension to Multiple Classes

RWTH AACHEN UNIVERSITY

- Two simple strategies
  - One-vs-all classifiers*
  - One-vs-one classifiers*

How many classifiers do we need in both cases?  
What difficulties do you see for those strategies?

B. Leibe 19  
Image source: C.M. Bishop, 2006

Machine Learning Winter '19

## Extension to Multiple Classes

RWTH AACHEN UNIVERSITY

- Problem
  - Both strategies result in regions 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 outputs  $y_k > y_j \forall j \neq k$ .
- Solution
  - We can avoid those difficulties by taking  $K$  linear functions of the form
 
$$y_k(\mathbf{x}) = \mathbf{w}_k^T \mathbf{x} + w_{k0}$$
 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
 
$$\mathbf{t}_k = (0, 1, 0, \dots, 0)^T$$

B. Leibe 20  
Image source: C.M. Bishop, 2006

Machine Learning Winter '19

## Extension to Multiple Classes

RWTH AACHEN UNIVERSITY

- 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**.

- Convex means: if  $\mathbf{x}_A$  and  $\mathbf{x}_B$  are both in  $\mathcal{R}_k$ , then any point  $\tilde{\mathbf{x}}$  on the connecting line between  $\mathbf{x}_A$  and  $\mathbf{x}_B$  is also in  $\mathcal{R}_k$ .

This makes linear discriminant models particularly suitable for problems for which the conditional densities  $p(\mathbf{x}|w_j)$  are unimodal.

B. Leibe 21  
Image source: C.M. Bishop, 2006

Machine Learning Winter '19

## Topics of This Lecture

RWTH AACHEN UNIVERSITY

- Linear discriminant functions
  - Definition
  - Extension to multiple classes
- Least-squares classification
  - Derivation
  - Shortcomings
- Generalized linear models
  - Connection to neural networks
  - Generalized linear discriminants & gradient descent

B. Leibe 22

Machine Learning Winter '19

## General Classification Problem

RWTH AACHEN UNIVERSITY

- 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
 
$$\mathbf{y}(\mathbf{x}) = \widetilde{\mathbf{W}}^T \widetilde{\mathbf{x}}$$
 where
 
$$\widetilde{\mathbf{W}} = [\widetilde{\mathbf{w}}_1, \dots, \widetilde{\mathbf{w}}_K] = \begin{bmatrix} w_{10} & \dots & w_{K0} \\ w_{11} & \dots & w_{K1} \\ \vdots & \ddots & \vdots \\ w_{1D} & \dots & w_{KD} \end{bmatrix}$$

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

B. Leibe 23

Machine Learning Winter '19

## General Classification Problem

RWTH AACHEN UNIVERSITY

- Classification problem
  - For the entire dataset, we can write
 
$$\mathbf{Y}(\widetilde{\mathbf{X}}) = \widetilde{\mathbf{X}} \widetilde{\mathbf{W}}$$
 and compare this to the target matrix  $\mathbf{T}$  where
 
$$\widetilde{\mathbf{W}} = [\widetilde{\mathbf{w}}_1, \dots, \widetilde{\mathbf{w}}_K]$$

$$\widetilde{\mathbf{X}} = \begin{bmatrix} \mathbf{x}_1^T \\ \vdots \\ \mathbf{x}_N^T \end{bmatrix} \quad \mathbf{T} = \begin{bmatrix} \mathbf{t}_1^T \\ \vdots \\ \mathbf{t}_N^T \end{bmatrix}$$

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

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

B. Leibe 24

RWTH AACHEN UNIVERSITY

## Least-Squares Classification

- Simplest approach
  - Directly try to minimize the **sum-of-squares error**
  - We could write this as

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

$$= \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (\mathbf{w}_k^T \mathbf{x}_n - t_{kn})^2$$

- 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...)

25

RWTH AACHEN UNIVERSITY

## Least-Squares Classification

- Multi-class case
  - Let's formulate the **sum-of-squares error** in matrix notation

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

- Taking the derivative yields

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

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

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

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

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

chain rule:  
 $\frac{\partial \mathbf{Z}}{\partial \mathbf{X}} = \frac{\partial \mathbf{Z}}{\partial \mathbf{Y}} \frac{\partial \mathbf{Y}}{\partial \mathbf{X}}$

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

28

RWTH AACHEN UNIVERSITY

## Least-Squares Classification

- Minimizing the sum-of-squares error

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

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

$$\tilde{\mathbf{W}} = (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}})^{-1} \tilde{\mathbf{X}}^T \mathbf{T}$$

$$= \tilde{\mathbf{X}}^\dagger \mathbf{T} \quad \text{"pseudo-inverse"}$$

- We then obtain the discriminant function as

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

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

29

RWTH AACHEN UNIVERSITY

## Problems with Least Squares

- Least-squares is very sensitive to outliers!
  - The error function penalizes predictions that are "too correct".

30

RWTH AACHEN UNIVERSITY

## 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!

31

RWTH AACHEN UNIVERSITY

## Topics of This Lecture

- Linear discriminant functions
  - Definition
  - Extension to multiple classes
- Least-squares classification
  - Derivation
  - Shortcomings
- Generalized linear models
  - Connection to neural networks
  - Generalized linear discriminants & gradient descent

32

RWTH AACHEN UNIVERSITY

## 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 may be nonlinear.
  - The decision surfaces 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  $\mathbf{x}$ .

33

RWTH AACHEN UNIVERSITY

## Generalized Linear Models

- Consider 2 classes:
 
$$p(\mathcal{C}_1 | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1) + p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)}$$

$$= \frac{1}{1 + \frac{p(\mathbf{x} | \mathcal{C}_2)p(\mathcal{C}_2)}{p(\mathbf{x} | \mathcal{C}_1)p(\mathcal{C}_1)}}$$

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

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

34

RWTH AACHEN UNIVERSITY

## Logistic Sigmoid Activation Function

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

Example: Normal distributions with identical covariance

The figure shows three plots. On the left is the logistic sigmoid function  $g(a) = 1 / (1 + \exp(-a))$ , which is an S-shaped curve ranging from 0 to 1. On the right are two plots. The top plot shows two normal distributions,  $p(x|a)$  (green) and  $p(x|b)$  (red), centered at different values of  $x$ . The bottom plot shows the corresponding conditional probabilities  $p(a|x)$  (green) and  $p(b|x)$  (red), which are sigmoidal curves that sum to 1.

35

RWTH AACHEN UNIVERSITY

## Normalized Exponential

- General case of  $K > 2$  classes:
 
$$p(\mathcal{C}_k | \mathbf{x}) = \frac{p(\mathbf{x} | \mathcal{C}_k)p(\mathcal{C}_k)}{\sum_j p(\mathbf{x} | \mathcal{C}_j)p(\mathcal{C}_j)}$$

$$= \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

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 multiclass generalization of the logistic sigmoid.

36

RWTH AACHEN UNIVERSITY

## Relationship to Neural Networks

- 2-Class case
 
$$y(\mathbf{x}) = g\left(\sum_{i=0}^D w_i x_i\right)$$
 with  $x_0 = 1$  constant
- Neural network ("single-layer perceptron")

The diagram shows a single-layer perceptron. It has a bias input  $x_0 = 1$  and  $D$  feature inputs  $x_1, \dots, x_D$ . Each input is multiplied by a weight  $w_0, w_1, \dots, w_D$ . The weighted inputs are summed and passed through an activation function  $g$  to produce the output  $y(x)$ .

37

RWTH AACHEN UNIVERSITY

## Relationship to Neural Networks

- Multi-class case
 
$$y_k(\mathbf{x}) = g\left(\sum_{i=0}^D w_{ki} x_i\right)$$
 with  $x_0 = 1$  constant
- Multi-class perceptron

The diagram shows a multi-class perceptron. It has a bias input  $x_0 = 1$  and  $D$  feature inputs  $x_1, \dots, x_D$ . Each input is multiplied by a weight  $w_{k0}, w_{k1}, \dots, w_{ki}$  for each class  $k$ . The weighted inputs are summed and passed through an activation function  $g$  to produce the output  $y_k(x)$ .

38

RWTH AACHEN UNIVERSITY

## Logistic Discrimination

- If we use the logistic sigmoid activation function...
 
$$g(a) \equiv \frac{1}{1 + \exp(-a)}$$

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

threshold

$x_0 = 1$

weights

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

Machine Learning Winter '19 | Slide adapted from Bernt Schiele | B. Leibe | 39

RWTH AACHEN UNIVERSITY

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

Machine Learning Winter '19 | B. Leibe | 40

RWTH AACHEN UNIVERSITY

## 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.
- Disadvantage
  - Least-squares minimization in general no longer leads to a closed-form analytical solution.
    - ⇒ Need to apply iterative methods.
    - ⇒ Gradient descent.

Machine Learning Winter '19 | B. Leibe | 41

RWTH AACHEN UNIVERSITY

## Linear Separability

- Up to now: restrictive assumption
  - Only consider linear decision boundaries
- Classical counterexample: XOR

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 42

RWTH AACHEN UNIVERSITY

## Generalized Linear Discriminants

- Generalization
  - Transform vector  $\mathbf{x}$  with  $M$  nonlinear basis functions  $\phi_j(\mathbf{x})$ :
 
$$y_k(\mathbf{x}) = \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}) + w_{k0}$$
  - Purpose of  $\phi_j(\mathbf{x})$ : basis functions
  - Allow non-linear decision boundaries.
  - By choosing the right  $\phi_j$ , every continuous function can (in principle) be approximated with arbitrary accuracy.
- Notation
 
$$y_k(\mathbf{x}) = \sum_{j=0}^M w_{kj} \phi_j(\mathbf{x}) \quad \text{with} \quad \phi_0(\mathbf{x}) = 1$$

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 44

RWTH AACHEN UNIVERSITY

## Generalized Linear Discriminants

- Model
 
$$y_k(\mathbf{x}) = \sum_{j=0}^M w_{kj} \phi_j(\mathbf{x}) = y_k(\mathbf{x}; \mathbf{w})$$
  - $K$  functions (outputs)  $y_k(\mathbf{x}; \mathbf{w})$
- Learning in Neural Networks
  - Single-layer networks:  $\phi_j$  are fixed, only weights  $\mathbf{w}$  are learned.
  - Multi-layer networks: both the  $\mathbf{w}$  and the  $\phi_j$  are learned.

... We will take a closer look at neural networks from lecture 11 on. For now, let's first consider generalized linear discriminants in general...

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 45

RWTH AACHEN UNIVERSITY

## Gradient Descent

- Learning the weights  $\mathbf{w}$ :
  - $N$  training data points:  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
  - $K$  outputs of decision functions:  $y_k(\mathbf{x}_n; \mathbf{w})$
  - Target vector for each data point:  $\mathbf{T} = \{t_1, \dots, t_N\}$
- Error function (least-squares error) of linear model
 
$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn})^2$$

$$= \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left( \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}_n) - t_{kn} \right)^2$$

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 46

RWTH AACHEN UNIVERSITY

## Gradient Descent

- Problem
  - The error function can in general no longer be minimized in closed form.
- Idea (Gradient Descent)
  - Iterative minimization
  - Start with an initial guess for the parameter values  $w_{kj}^{(0)}$
  - Move towards a (local) minimum by following the gradient.
 
$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

$\eta$ : Learning rate
  - This simple scheme corresponds to a 1<sup>st</sup>-order Taylor expansion (There are more complex procedures available).

Machine Learning Winter '19 | B. Leibe | 47

RWTH AACHEN UNIVERSITY

## Gradient Descent – Basic Strategies

- “Batch learning”
 
$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

$\eta$ : Learning rate

  - Compute the gradient based on all training data:
 
$$\frac{\partial E(\mathbf{w})}{\partial w_{kj}}$$

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 48

RWTH AACHEN UNIVERSITY

## Gradient Descent – Basic Strategies

- “Sequential updating”
 
$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \left. \frac{\partial E_n(\mathbf{w})}{\partial w_{kj}} \right|_{\mathbf{w}^{(\tau)}}$$

$\eta$ : Learning rate

  - Compute the gradient based on a single data point at a time:
 
$$\frac{\partial E_n(\mathbf{w})}{\partial w_{kj}}$$

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 49

RWTH AACHEN UNIVERSITY

## Gradient Descent

- Error function
 
$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \sum_{k=1}^K \left( \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}_n) - t_{kn} \right)^2$$

$$E_n(\mathbf{w}) = \frac{1}{2} \sum_{k=1}^K \left( \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}_n) - t_{kn} \right)^2$$

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{kj}} = \left( \sum_{j=1}^M w_{kj} \phi_j(\mathbf{x}_n) - t_{kn} \right) \phi_j(\mathbf{x}_n)$$

$$= (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn}) \phi_j(\mathbf{x}_n)$$

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 50

RWTH AACHEN UNIVERSITY

## Gradient Descent

- Delta rule (=LMS rule)
 
$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn}) \phi_j(\mathbf{x}_n)$$

$$= w_{kj}^{(\tau)} - \eta \delta_{kn} \phi_j(\mathbf{x}_n)$$
  - where
 
$$\delta_{kn} = y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn}$$
  - $\Rightarrow$  Simply feed back the input data point, weighted by the classification error.

Machine Learning Winter '19 | Slide credit: Bernt Schiele | B. Leibe | 51



## Gradient Descent

- Cases with differentiable, non-linear activation function

$$y_k(\mathbf{x}) = g(a_k) = g\left(\sum_{j=0}^M w_{kj} \phi_j(\mathbf{x}_n)\right)$$

- Gradient descent

$$\frac{\partial E_n(\mathbf{w})}{\partial w_{kj}} = \frac{\partial g(a_k)}{\partial w_{kj}} (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn}) \phi_j(\mathbf{x}_n)$$

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta \delta_{kn} \phi_j(\mathbf{x}_n)$$

$$\delta_{kn} = \frac{\partial g(a_k)}{\partial w_{kj}} (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn})$$

## Summary: Generalized Linear Discriminants

- Properties
  - General class of decision functions.
  - Nonlinearity  $g(\cdot)$  and basis functions  $\phi_j$  allow us to address linearly non-separable problems.
  - Shown simple sequential learning approach for parameter estimation using gradient descent.
  - Better 2<sup>nd</sup> order gradient descent approaches available (e.g. Newton-Raphson).
- Limitations / Caveats
  - Flexibility of model is limited by curse of dimensionality
    - $g(\cdot)$  and  $\phi_j$  often introduce additional parameters.
    - Models are either limited to lower-dimensional input space or need to share parameters.
  - Linearly separable case often leads to overfitting.
    - Several possible parameter choices minimize training error.

## References and Further Reading

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

Christopher M. Bishop  
Pattern Recognition and Machine Learning  
Springer, 2006

