

RWTH AACHEN
UNIVERSITY

Machine Learning - Lecture 6

Linear Discriminants II

09.05.2016

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

Machine Learning, Summer '16

RWTH AACHEN
UNIVERSITY

Course Outline

- Fundamentals (2 weeks)
 - Bayes Decision Theory
 - Probability Density Estimation
- Discriminative Approaches (5 weeks)
 - Linear Discriminant Functions
 - Support Vector Machines
 - Ensemble Methods & Boosting
 - Randomized Trees, Forests & Ferns
- Generative Models (4 weeks)
 - Bayesian Networks
 - Markov Random Fields

B. Leibe

2

Machine Learning, Summer '16

RWTH AACHEN
UNIVERSITY

Recap: Linear Discriminant Functions

- Basic idea
 - Directly encode decision boundary
 - Minimize misclassification probability directly.
- Linear discriminant functions
 - $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - w , w_0 define a hyperplane in \mathbb{R}^D .
 - If a data set can be perfectly classified by a linear discriminant, then we call it **linearly separable**.

B. Leibe

3

Machine Learning, Summer '16

Slide adapted from Bernt Schiele

RWTH AACHEN
UNIVERSITY

Recap: Least-Squares Classification

- Simplest approach
 - Directly try to minimize the **sum-of-squares error**

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

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

- Setting the derivative to zero yields

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

- We then obtain the discriminant function as

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

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

B. Leibe

4

Machine Learning, Summer '16

RWTH AACHEN
UNIVERSITY

Recap: Problems with Least Squares

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

B. Leibe

5

Machine Learning, Summer '16

Image source: C. M. Bishop, 2006

RWTH AACHEN
UNIVERSITY

Recap: Generalized Linear Models

- 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} .
- Advantages of the non-linearity
 - Can be used to bound the influence of outliers and “too correct” data points.
 - When using a sigmoid for $g(\cdot)$, we can interpret the $y(\mathbf{x})$ as posterior probabilities.

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

B. Leibe

6

Machine Learning, Summer '16

RWTH AACHEN UNIVERSITY

Recap: Linear Separability

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

7

Machine Learning, Summer '16
Slide credit: Bernt Schiele
B. Leibe

RWTH AACHEN UNIVERSITY

Linear Separability

- Even if the data is not linearly separable, a linear decision boundary may still be "optimal".
 - Generalization
 - E.g. in the case of Normal distributed data (with equal covariance matrices)
- Choice of the right discriminant function is important and should be based on
 - Prior knowledge (of the general functional form)
 - Empirical comparison of alternative models
 - Linear discriminants are often used as benchmark.

9

Machine Learning, Summer '16
Slide credit: Bernt Schiele
B. Leibe

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 ϕ_j , every continuous function can (in principle) be approximated with arbitrary accuracy.
- Notation

$$y_k(\mathbf{x}) = \sum_{j=0}^M w_{kj} \phi_j(\mathbf{x}) \quad \text{with } \phi_0(\mathbf{x}) = 1$$

10

Machine Learning, Summer '16
Slide credit: Bernt Schiele
B. Leibe

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: ϕ_j are fixed, only weights \mathbf{w} are learned.
 - Multi-layer networks: both the \mathbf{w} and the ϕ_j are learned.
 - In the following, we will not go into details about neural networks in particular, but consider generalized linear discriminants in general...

11

Machine Learning, Summer '16
Slide credit: Bernt Schiele
B. Leibe

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

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

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

12

Machine Learning, Summer '16
Slide credit: Bernt Schiele
B. Leibe

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.

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

η : Learning rate
 - This simple scheme corresponds to a 1st-order Taylor expansion (There are more complex procedures available).

13

Machine Learning, Summer '16
B. Leibe

RWTH AACHEN UNIVERSITY

Gradient Descent - Basic Strategies

- “Batch learning”

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

η : Learning rate

- Compute the gradient based on all training data:

$$\frac{\partial E(\mathbf{w})}{\partial w_{kj}}$$

Machine Learning, Summer '16 17
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Gradient Descent - Basic Strategies

- “Sequential updating”

$$E(\mathbf{w}) = \sum_{n=1}^N E_n(\mathbf{w})$$

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

η : 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, Summer '16 18
Slide credit: Bernt Schiele B. Leibe

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, Summer '16 19
Slide credit: Bernt Schiele B. Leibe

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

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

⇒ Simply feed back the input data point, weighted by the classification error.

Machine Learning, Summer '16 20
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

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

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

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

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

Machine Learning, Summer '16 21
Slide credit: Bernt Schiele B. Leibe

RWTH AACHEN UNIVERSITY

Summary: Generalized Linear Discriminants

- Properties
 - General class of decision functions.
 - Nonlinearity $g(\cdot)$ and basis functions ϕ_j allow us to address linearly non-separable problems.
 - Shown simple sequential learning approach for parameter estimation using gradient descent.
 - Better 2nd order gradient descent approaches available (e.g. Newton-Raphson).
- Limitations / Caveats
 - Flexibility of model is limited by curse of dimensionality
 - $g(\cdot)$ and ϕ_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.

Machine Learning, Summer '16 22

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Fisher's linear discriminant (FLD)
 - Classification as dimensionality reduction
 - Linear discriminant analysis
 - Multiple discriminant analysis
 - Applications
- Logistic Regression
 - Probabilistic discriminative models
 - Logistic sigmoid (logit function)
 - Cross-entropy error
 - Gradient descent
 - Iteratively Reweighted Least Squares
- Note on Error Functions

Machine Learning, Summer '16 23

RWTH AACHEN UNIVERSITY

Classification as Dimensionality Reduction

- Classification as dimensionality reduction
 - We can interpret the linear classification model as a projection onto a lower-dimensional space.
 - E.g., take the D -dimensional input vector x and project it down to one dimension by applying the function

$$y = w^T x$$
 - If we now place a threshold at $y \geq -w_0$, we obtain our standard two-class linear classifier.
 - The classifier will have a lower error the better this projection separates the two classes.

⇒ New interpretation of the learning problem

- Try to find the projection vector w that maximizes the class separation.

Machine Learning, Summer '16 24

RWTH AACHEN UNIVERSITY

Classification as Dimensionality Reduction

- Two questions
 - How to measure class separation?
 - How to find the best projection (with maximal class separation)?

Machine Learning, Summer '16 25

RWTH AACHEN UNIVERSITY

Classification as Dimensionality Reduction

- Measuring class separation
 - We could simply measure the separation of the class means.
 - ⇒ Choose w so as to maximize

$$(m_2 - m_1) = w^T (m_2 - m_1)$$
- Problems with this approach
 1. This expression can be made arbitrarily large by increasing $\|w\|$.
 ⇒ Need to enforce additional constraint $\|w\| = 1$.
 ⇒ This constrained minimization results in $w \propto (m_2 - m_1)$
 2. The criterion may result in bad separation if the clusters have elongated shapes.

Machine Learning, Summer '16 26

RWTH AACHEN UNIVERSITY

Fisher's Linear Discriminant Analysis (FLD)

- Better idea:
 - Find a projection that maximizes the ratio of the between-class variance to the within-class variance:

$$J(w) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2} \quad \text{with} \quad s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$$
 - Usually, this is written as

$$J(w) = \frac{w^T S_B w}{w^T S_W w}$$
 - where

$$S_B = (m_2 - m_1)(m_2 - m_1)^T$$
 between-class scatter matrix
 - $$S_W = \sum_{k=1}^2 \sum_{n \in C_k} (x_n - m_k)(x_n - m_k)^T$$
 within-class scatter matrix

Machine Learning, Summer '16 27

RWTH AACHEN UNIVERSITY

Fisher's Linear Discriminant Analysis (FLD)

- Maximize distance between classes
- Minimize distance within a class
- Criterion: $J(w) = \frac{w^T S_B w}{w^T S_W w}$
- S_B ... between-class scatter matrix
 S_W ... within-class scatter matrix
- The optimal solution for w can be obtained as:

$$w \propto S_W^{-1} (m_2 - m_1)$$
- Classification function:

$$y(x) = w^T x + w_0 \begin{cases} \text{Class 1} & \geq 0 \\ \text{Class 2} & < 0 \end{cases}$$

Machine Learning, Summer '16 28

Slide adapted from Ales Leonardis

Machine Learning, Summer '16

Multiple Discriminant Analysis

- Generalization to K classes

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$

where

$$\mathbf{W} = [\mathbf{w}_1, \dots, \mathbf{w}_K] \quad \mathbf{m} = \frac{1}{N} \sum_{n=1}^N \mathbf{x}_n = \frac{1}{N} \sum_{k=1}^K N_k \mathbf{m}_k$$

$$\mathbf{S}_B = \sum_{k=1}^K N_k (\mathbf{m}_k - \mathbf{m})(\mathbf{m}_k - \mathbf{m})^T$$

$$\mathbf{S}_W = \sum_{k=1}^K \sum_{n \in C_k} (\mathbf{x}_n - \mathbf{m}_k)(\mathbf{x}_n - \mathbf{m}_k)^T$$

B. Leibe 29

Machine Learning, Summer '16

Maximizing $J(\mathbf{W})$

- "Rayleigh quotient" \Rightarrow Generalized eigenvalue problem

$$J(\mathbf{W}) = \frac{|\mathbf{W}^T \mathbf{S}_B \mathbf{W}|}{|\mathbf{W}^T \mathbf{S}_W \mathbf{W}|}$$

- The columns of the optimal \mathbf{W} are the eigenvectors corresponding to the largest eigenvalues of $\mathbf{S}_B \mathbf{w}_i = \lambda_i \mathbf{S}_W \mathbf{w}_i$
- Defining $\mathbf{v} = \mathbf{S}_B^{-\frac{1}{2}} \mathbf{w}$, we get $\mathbf{S}_B^{-\frac{1}{2}} \mathbf{S}_W^{-1} \mathbf{S}_B^{-\frac{1}{2}} \mathbf{v} = \lambda \mathbf{v}$ which is a regular eigenvalue problem. \Rightarrow Solve to get eigenvectors of \mathbf{v} , then from that of \mathbf{w} .
- For the K -class case we obtain (at most) $K-1$ projections. (i.e. eigenvectors corresponding to non-zero eigenvalues.)

B. Leibe 30

Machine Learning, Summer '16


What Does It Mean?

- What does it mean to apply a linear classifier?

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

Weight vector Input vector

- Classifier interpretation
 - The weight vector has the same dimensionality as \mathbf{x} .
 - Positive contributions where $\text{sign}(x_i) = \text{sign}(w_i)$.
 - \Rightarrow The weight vector identifies which input dimensions are important for positive or negative classification (large $|w_i|$) and which ones are irrelevant (near-zero w_i).
 - \Rightarrow If the inputs \mathbf{x} are normalized, we can interpret \mathbf{w} as a "template" vector that the classifier tries to match.

$$\mathbf{w}^T \mathbf{x} = \|\mathbf{w}\| \|\mathbf{x}\| \cos \theta$$


B. Leibe 31

Machine Learning, Summer '16

Example Application: Fisherfaces

- Visual discrimination task
 - Training data:
 - C_1 : Subjects with glasses
 - C_2 : Subjects without glasses
 - Test:
 - Image of a person with glasses: glasses?
 - Diagram showing an image being converted to a vector of pixel values.
 - Text: "Take each image as a vector of pixel values and apply FLD..."


B. Leibe 32
Image source: Yale Face Database

Machine Learning, Summer '16


Fisherfaces: Interpretability

- Resulting weight vector for "Glasses/NoGlasses"

\mathbf{x}



\mathbf{w}



Slide credit: Peter Belhumeur B. Leibe [Belhumeur et al., 1997]

Machine Learning, Summer '16

Summary: Fisher's Linear Discriminant

- Properties
 - Simple method for dimensionality reduction, preserves class discriminability.
 - Can use parametric methods in reduced-dim. space that might not be feasible in original higher-dim. space.
 - Widely used in practical applications.
- Restrictions / Caveats
 - Not possible to get more than $K-1$ projections.
 - FLD reduces the computation to class means and covariances. \Rightarrow Implicit assumption that class distributions are unimodal and well-approximated by a Gaussian/hyperellipsoid.

B. Leibe 34

RWTH AACHEN
UNIVERSITY

Topics of This Lecture

- Fisher's linear discriminant (FLD)
 - Classification as dimensionality reduction
 - Linear discriminant analysis
 - Multiple discriminant analysis
 - Applications
- Logistic Regression
 - Probabilistic discriminative models
 - Logistic sigmoid (logit function)
 - Cross-entropy error
 - Gradient descent
 - Iteratively Reweighted Least Squares
- Note on Error Functions

Machine Learning, Summer '16 35

RWTH AACHEN
UNIVERSITY

Probabilistic Discriminative Models

- We have seen that we can write

$$p(\mathcal{C}_1|\mathbf{x}) = \sigma(a) \quad \text{logistic sigmoid function}$$

$$= \frac{1}{1 + \exp(-a)}$$
- We can obtain the familiar probabilistic model by setting

$$a = \ln \frac{p(\mathbf{x}|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\mathbf{x}|\mathcal{C}_2)p(\mathcal{C}_2)}$$
- Or we can use generalized linear discriminant models

$$a = \mathbf{w}^T \mathbf{x}$$
 or

$$a = \mathbf{w}^T \phi(\mathbf{x})$$

Machine Learning, Summer '16 36

RWTH AACHEN
UNIVERSITY

Probabilistic Discriminative Models

- In the following, we will consider models of the form

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$
 with

$$p(\mathcal{C}_2|\phi) = 1 - p(\mathcal{C}_1|\phi)$$
- This model is called **logistic regression**.
- Why should we do this? What advantage does such a model have compared to modeling the probabilities?

$$p(\mathcal{C}_1|\phi) = \frac{p(\phi|\mathcal{C}_1)p(\mathcal{C}_1)}{p(\phi|\mathcal{C}_1)p(\mathcal{C}_1) + p(\phi|\mathcal{C}_2)p(\mathcal{C}_2)}$$
- Any ideas?

Machine Learning, Summer '16 37

RWTH AACHEN
UNIVERSITY

Comparison

- Let's look at the number of parameters...
 - Assume we have an M -dimensional feature space ϕ .
 - And assume we represent $p(\phi|\mathcal{C}_k)$ and $p(\mathcal{C}_k)$ by Gaussians.
 - How many parameters do we need?
 - For the means: $2M$
 - For the covariances: $M(M+1)/2$
 - Together with the class priors, this gives $M(M+5)/2+1$ parameters!
 - How many parameters do we need for logistic regression?

$$p(\mathcal{C}_1|\phi) = y(\phi) = \sigma(\mathbf{w}^T \phi)$$
 - Just the values of $\mathbf{w} \Rightarrow M$ parameters.

\Rightarrow For large M , logistic regression has clear advantages!

Machine Learning, Summer '16 38

RWTH AACHEN
UNIVERSITY

Logistic Sigmoid

- Properties
 - Definition: $\sigma(a) = \frac{1}{1 + \exp(-a)}$
 - Inverse: $a = \ln \left(\frac{\sigma}{1 - \sigma} \right)$ "logit" function
 - Symmetry property:

$$\sigma(-a) = 1 - \sigma(a)$$
 - Derivative: $\frac{d\sigma}{da} = \sigma(1 - \sigma)$

Machine Learning, Summer '16 39

RWTH AACHEN
UNIVERSITY

Logistic Regression

- Let's consider a data set $\{\phi_n, t_n\}$ with $n = 1, \dots, N$, where $\phi_n = \phi(\mathbf{x}_n)$ and $t_n \in \{0, 1\}$, $\mathbf{t} = (t_1, \dots, t_N)^T$.
- With $y_n = p(\mathcal{C}_1|\phi_n)$, we can write the likelihood as

$$p(\mathbf{t}|\mathbf{w}) = \prod_{n=1}^N y_n^{t_n} (1 - y_n)^{1-t_n}$$
- Define the error function as the negative log-likelihood

$$E(\mathbf{w}) = -\ln p(\mathbf{t}|\mathbf{w})$$

$$= -\sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$
 - This is the so-called **cross-entropy error function**.

Machine Learning, Summer '16 40

RWTH AACHEN UNIVERSITY

Gradient of the Error Function

$$y_n = \sigma(\mathbf{w}^T \phi_n)$$

$$\frac{dy_n}{d\mathbf{w}} = y_n(1 - y_n)\phi_n$$

- Error function**

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$
- Gradient**

$$\begin{aligned} \nabla E(\mathbf{w}) &= - \sum_{n=1}^N \left\{ t_n \frac{d}{d\mathbf{w}} y_n + (1 - t_n) \frac{d}{d\mathbf{w}} (1 - y_n) \right\} \\ &= - \sum_{n=1}^N \left\{ t_n \frac{dy_n}{d\mathbf{w}} - (1 - t_n) \frac{dy_n}{d\mathbf{w}} \right\} \\ &= - \sum_{n=1}^N \{ (t_n - (1 - t_n)) y_n \phi_n \} \\ &= \sum_{n=1}^N (y_n - t_n) \phi_n \end{aligned}$$

41

RWTH AACHEN UNIVERSITY

Gradient of the Error Function

- Gradient for logistic regression**

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n$$
- Does this look familiar to you?
- This is the same result as for the Delta (=LMS) rule

$$w_{kj}^{(\tau+1)} = w_{kj}^{(\tau)} - \eta (y_k(\mathbf{x}_n; \mathbf{w}) - t_{kn}) \phi_j(\mathbf{x}_n)$$
- We can use this to derive a sequential estimation algorithm.
 - However, this will be quite slow...

42

RWTH AACHEN UNIVERSITY

A More Efficient Iterative Method...

- Second-order Newton-Raphson gradient descent scheme**

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \mathbf{H}^{-1} \nabla E(\mathbf{w})$$

where $\mathbf{H} = \nabla \nabla E(\mathbf{w})$ is the Hessian matrix, i.e. the matrix of second derivatives.
- Properties**
 - Local quadratic approximation to the log-likelihood.
 - Faster convergence.

43

RWTH AACHEN UNIVERSITY

Newton-Raphson for Least-Squares Estimation

- Let's first apply Newton-Raphson to the least-squares error function:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n)^2$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (\mathbf{w}^T \phi_n - t_n) \phi_n = \Phi^T \Phi \mathbf{w} - \Phi^T \mathbf{t}$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N \phi_n \phi_n^T = \Phi^T \Phi \quad \text{where } \Phi = \begin{bmatrix} \phi_1^T \\ \vdots \\ \phi_N^T \end{bmatrix}$$
- Resulting update scheme:**

$$\begin{aligned} \mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - (\Phi^T \Phi)^{-1} (\Phi^T \Phi \mathbf{w}^{(\tau)} - \Phi^T \mathbf{t}) \\ &= (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \end{aligned}$$

Closed-form solution!

44

RWTH AACHEN UNIVERSITY

Newton-Raphson for Logistic Regression

- Now, let's try Newton-Raphson on the cross-entropy error function:

$$E(\mathbf{w}) = - \sum_{n=1}^N \{t_n \ln y_n + (1 - t_n) \ln(1 - y_n)\}$$

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n = \Phi^T (\mathbf{y} - \mathbf{t})$$

$$\mathbf{H} = \nabla \nabla E(\mathbf{w}) = \sum_{n=1}^N y_n(1 - y_n) \phi_n \phi_n^T = \Phi^T \mathbf{R} \Phi$$

where \mathbf{R} is an $N \times N$ diagonal matrix with $R_{nn} = y_n(1 - y_n)$.

⇒ The Hessian is no longer constant, but depends on \mathbf{w} through the weighting matrix \mathbf{R} .

45

RWTH AACHEN UNIVERSITY

Iteratively Reweighted Least Squares

- Update equations**

$$\begin{aligned} \mathbf{w}^{(\tau+1)} &= \mathbf{w}^{(\tau)} - (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T (\mathbf{y} - \mathbf{t}) \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \left\{ \Phi^T \mathbf{R} \Phi \mathbf{w}^{(\tau)} - \Phi^T (\mathbf{y} - \mathbf{t}) \right\} \\ &= (\Phi^T \mathbf{R} \Phi)^{-1} \Phi^T \mathbf{R} \mathbf{z} \end{aligned}$$

with $\mathbf{z} = \Phi \mathbf{w}^{(\tau)} - \mathbf{R}^{-1} (\mathbf{y} - \mathbf{t})$
- Again very similar form (normal equations)**
 - But now with non-constant weighting matrix \mathbf{R} (depends on \mathbf{w}).
 - Need to apply normal equations iteratively.
 - ⇒ Iteratively Reweighted Least-Squares (IRLS)

46

RWTH AACHEN UNIVERSITY

Summary: Logistic Regression

- **Properties**
 - Directly represent posterior distribution $p(\phi | C_k)$
 - Requires fewer parameters than modeling the likelihood + prior.
 - Very often used in statistics.
 - It can be shown that the cross-entropy error function is concave
 - Optimization leads to unique minimum
 - But no closed-form solution exists
 - Iterative optimization (IRLS)
 - Both online and batch optimizations exist
 - There is a multi-class version described in (Bishop Ch.4.3.4).
- **Caveat**
 - Logistic regression tends to systematically overestimate odds ratios when the sample size is less than -500.

47

RWTH AACHEN UNIVERSITY

Topics of This Lecture

- Fisher's linear discriminant (FLD)
 - Classification as dimensionality reduction
 - Linear discriminant analysis
 - Multiple discriminant analysis
 - Applications
- Logistic Regression
 - Probabilistic discriminative models
 - Logistic sigmoid (logit function)
 - Cross-entropy error
 - Gradient descent
 - Iteratively Reweighted Least Squares
- **Note on Error Functions**

48

RWTH AACHEN UNIVERSITY

Note on Error Functions

$t_n \in \{ -1, 1 \}$

Ideal misclassification error

- **Ideal misclassification error function (black)**
 - This is what we want to approximate,
 - Unfortunately, it is not differentiable.
 - The gradient is zero for misclassified points.
 - ⇒ We cannot minimize it by gradient descent.

49

RWTH AACHEN UNIVERSITY

Note on Error Functions

$t_n \in \{ -1, 1 \}$

Ideal misclassification error
Squared error

- **Squared error used in Least-Squares Classification**
 - Very popular, leads to closed-form solutions.
 - However, sensitive to outliers due to squared penalty.
 - Penalizes "too correct" data points
 - ⇒ Generally does not lead to good classifiers.

50

RWTH AACHEN UNIVERSITY

Comparing Error Functions (Loss Functions)

$t_n \in \{ -1, 1 \}$

Ideal misclassification error
Squared error
Cross-entropy error

- **Cross-Entropy Error**
 - Minimizer of this error is given by posterior class probabilities.
 - Concave error function, unique minimum exists.
 - Robust to outliers, error increases only roughly linearly
 - But no closed-form solution, requires iterative estimation.

51

RWTH AACHEN UNIVERSITY

Overview: Error Functions

- **Ideal Misclassification Error**
 - This is what we would like to optimize.
 - But cannot compute gradients here.
- **Quadratic Error**
 - Easy to optimize, closed-form solutions exist.
 - But not robust to outliers.
- **Cross-Entropy Error**
 - Minimizer of this error is given by posterior class probabilities.
 - Concave error function, unique minimum exists.
 - But no closed-form solution, requires iterative estimation.

⇒ **Analysis tool to compare classification approaches**

52

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 - 4.3).

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

