

Cognitive Algorithms

Lecture 3

Linear Classification

Klaus-Robert Müller
Felix Bießmann, Irene Winkler

Berlin Institute of Technology
Dept. Machine Learning

Recap

LDA

Bayes View

BBCI

Cross-validation

Summary

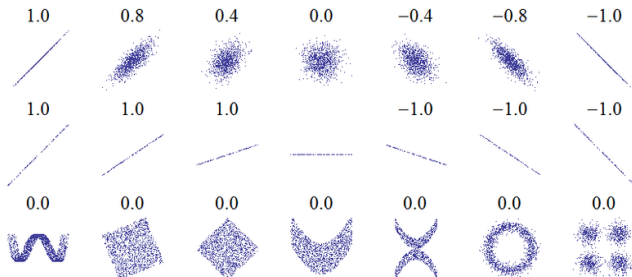
Covariance and Correlation

For two random variables X and Y , their **covariance** and **correlation** are defined as

$$\text{Cov}(X, Y) := E[(X - E(X))(Y - E(Y))]$$

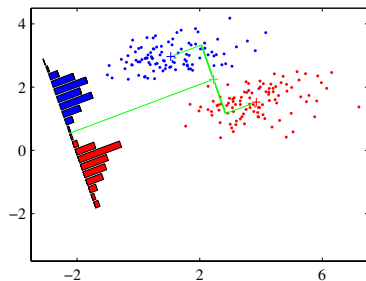
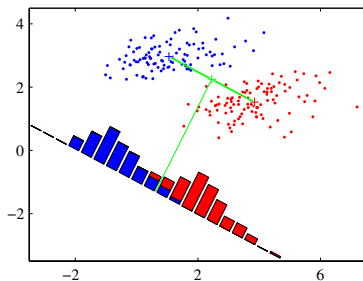
$$\text{Corr}(X, Y) := \frac{\text{Cov}(X, Y)}{\sqrt{\text{Var}(X)}\sqrt{\text{Var}(Y)}}.$$

Correlation measures the linear relationship between X and Y:



Linear Discriminant Analysis

View classification in terms of dimensionality reduction



Goal: Find a (normal vector of a linear decision boundary) \mathbf{w} that

- Maximizes mean class difference, and
- Minimizes variance in each class

Linear Discriminant Analysis

Goal: Find a (normal vector of a linear decision boundary) \mathbf{w} that
Maximizes mean class difference

$$(\mathbf{w}^\top \mathbf{w}_o - \mathbf{w}^\top \mathbf{w}_\Delta)^2 = \mathbf{w}^\top \underbrace{(\mathbf{w}_o - \mathbf{w}_\Delta)(\mathbf{w}_o - \mathbf{w}_\Delta)^\top}_{S_B - \text{"between class scatter"}} \mathbf{w} \quad (2)$$

Minimizes variance in each class

$$\begin{aligned} & \frac{1}{N_o} \sum_{i=1}^{N_o} \left(\mathbf{w}^\top (\mathbf{x}_{oi} - \mathbf{w}_o) \right)^2 + \frac{1}{N_\Delta} \sum_{j=1}^{N_\Delta} \left(\mathbf{w}^\top (\mathbf{x}_{\Delta j} - \mathbf{w}_\Delta) \right)^2 \\ &= \mathbf{w}^\top \underbrace{\left(\frac{1}{N_o} \sum_{i=1}^{N_o} (\mathbf{x}_{oi} - \mathbf{w}_o)(\mathbf{x}_{oi} - \mathbf{w}_o)^\top + \frac{1}{N_\Delta} \sum_{j=1}^{N_\Delta} (\mathbf{x}_{\Delta j} - \mathbf{w}_{\Delta j})(\mathbf{x}_{\Delta j} - \mathbf{w}_{\Delta j})^\top \right)}_{S_W - \text{"within class scatter"}} \mathbf{w} \end{aligned}$$

Linear Discriminant Analysis

Goal: Find a (normal vector of a linear decision boundary) \mathbf{w} that

Maximizes mean class difference, $\mathbf{w}^\top S_B \mathbf{w}$ and

Minimizes variance in each class, $\mathbf{w}^\top S_W \mathbf{w}$

→ maximize the *Fisher criterion*

$$\operatorname{argmax}_{\mathbf{w}} \frac{\mathbf{w}^\top S_B \mathbf{w}}{\mathbf{w}^\top S_W \mathbf{w}} \quad (3)$$

Linear Discriminant Analysis

$$\operatorname{argmax}_{\mathbf{w}} \frac{\mathbf{w}^{\top} S_B \mathbf{w}}{\mathbf{w}^{\top} S_W \mathbf{w}}$$

To optimize the Fisher criterion, we set its derivative w.r.t \mathbf{w} to 0

$$\begin{aligned} \frac{(\mathbf{w}^{\top} S_W \mathbf{w}) S_B \mathbf{w} - (\mathbf{w}^{\top} S_B \mathbf{w}) S_W \mathbf{w}}{(\mathbf{w}^{\top} S_W \mathbf{w})^2} &= 0 \\ (\mathbf{w}^{\top} S_B \mathbf{w}) S_W \mathbf{w} &= (\mathbf{w}^{\top} S_W \mathbf{w}) S_B \mathbf{w} \\ S_W \mathbf{w} &= S_B \mathbf{w} \underbrace{\frac{\mathbf{w}^{\top} S_W \mathbf{w}}{\mathbf{w}^{\top} S_B \mathbf{w}}}_{\text{scalar}} \end{aligned}$$

Linear Discriminant Analysis

$$\begin{aligned} \operatorname{argmax}_{\mathbf{w}} \frac{\mathbf{w}^\top S_B \mathbf{w}}{\mathbf{w}^\top S_W \mathbf{w}} \\ \rightarrow S_W \mathbf{w} = S_B \mathbf{w} \lambda \end{aligned}$$

Note that

$$S_B \mathbf{w} = (\mathbf{w}_o - \mathbf{w}_\Delta) \underbrace{(\mathbf{w}_o - \mathbf{w}_\Delta)^\top \mathbf{w}}_{\text{scalar}}$$

thus left multiplying with S_W^{-1} yields

$$\mathbf{w} \propto S_W^{-1} (\mathbf{w}_o - \mathbf{w}_\Delta).$$

Linear Discriminant Analysis

An equivalent formulation of LDA is given by proving for the total covariance of the data $S = S_W + \frac{N_\Delta N_o}{(N_\Delta + N_o)} S_B$. Then:

$$S_W \mathbf{w} \propto S_B \mathbf{w}$$

$$\left(S - \frac{N_\Delta N_o}{N_\Delta + N_o} S_B\right) \mathbf{w} \propto S_B \mathbf{w}$$

$$S \mathbf{w} \propto S_B \mathbf{w}$$

$$\mathbf{w} \propto S^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta).$$

→ LDA first *decorrelates* the data
followed by nearest centroid classification

Linear Discriminant Analysis

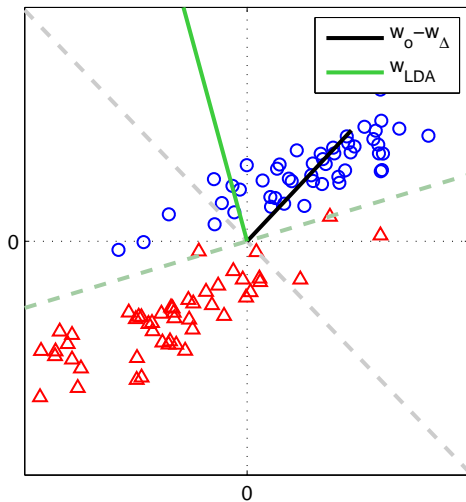
LDA first *decorrelates* the data followed by nearest centroid classification:

$$\begin{aligned}\mathbf{x} &\mapsto \text{sign}(\mathbf{w}^T \cdot \mathbf{x} - \beta) \\ \mathbf{w} &\propto S^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta)\end{aligned}$$

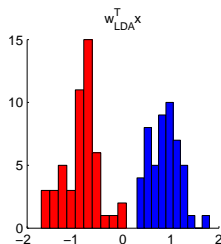
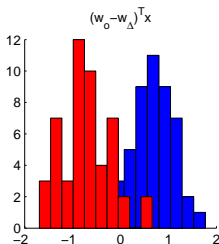
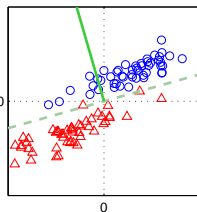
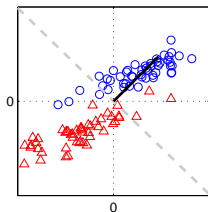
$$\mathbf{w}^T \mathbf{x} = (\mathbf{w}_o - \mathbf{w}_\Delta)^T S^{-1} \mathbf{x} = \underbrace{(\mathbf{w}_o - \mathbf{w}_\Delta)^T U \Lambda^{-1/2}}_{\text{mean class difference of decorrelated data}} \underbrace{\Lambda^{-1/2} U^T \mathbf{x}}_{\text{decorrelated } \mathbf{x}}$$

where $S = U \Lambda U^T$ is the eigenvalue decomposition of S

Linear Discriminant Analysis vs Nearest Centroid Classifier



Linear Discriminant Analysis vs Nearest Centroid Classifier



Bayesian decision theory

Bayesian decision theory:

For a new data point $\mathbf{x} \in \mathbb{R}^D$

Decide class Δ if $p(\Delta|\mathbf{x}) > p(o|\mathbf{x})$.

Calculate $p(\Delta|\mathbf{x})$ with Bayes rule:

$$\begin{aligned} p(\Delta|\mathbf{x}) &= \frac{p(\Delta, \mathbf{x})}{p(\mathbf{x})} \\ &= \frac{p(\Delta)p(\mathbf{x}|\Delta)}{p(\mathbf{x})} \end{aligned}$$

Bayesian decision theory

Estimating $p(\mathbf{x}|\Delta)$ is difficult: already if each dimension of \mathbf{x} can take 2 values $\rightarrow 2^D$ possible values.

One possibility to deal with it:

Choose a distribution $p(\mathbf{x}|\Delta)$, $p(\mathbf{x}|o)$ that is easy to deal with

\rightarrow Most popular: The Gaussian (or Normal) distribution

$$\mathbf{x} \in \mathbb{R}^D \sim \mathcal{N}(\mathbf{w}_\Delta, S_\Delta) = \frac{1}{(2\pi)^{\frac{D}{2}} \sqrt{|S_\Delta|}} e^{-\frac{1}{2}(\mathbf{x}-\mathbf{w}_\Delta)^\top S_\Delta^{-1}(\mathbf{x}-\mathbf{w}_\Delta)}$$

Linear Discriminant - a Probabilistic View

If we assume equal covariance in each class, $S_W = 2S_\Delta = 2S_o$, and equal class probabilities, $p(\Delta) = p(o) = 0.5$, the optimal classification boundary is linear and given by

$$\begin{aligned}\mathbf{w} &= S_W^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta) \\ \beta &= \frac{1}{2}\mathbf{w}_o S_W^{-1}\mathbf{w}_o + \frac{1}{2}\mathbf{w}_\Delta S_W^{-1}\mathbf{w}_\Delta = \frac{1}{2}\mathbf{w}^T(\mathbf{w}_o + \mathbf{w}_\Delta)\end{aligned}$$

⇒ Linear decision boundaries arise from simple assumption about the distribution of the data.

Linear Discriminant Algorithm

Computes: Normal vector \mathbf{w} of decision hyperplane, threshold β

Input: Data $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N)\}$, $\mathbf{x}_i \in \mathbb{R}^D$, $y_i \in \{-1, +1\}$,

 Compute class mean vectors

$$\mathbf{w}_{-1} = 1/N_- \sum_{i \in \mathcal{Y}_{-1}} \mathbf{x}_i$$

$$\mathbf{w}_{+1} = 1/N_+ \sum_{j \in \mathcal{Y}_{+1}} \mathbf{x}_j$$

 Compute *within-class* covariance matrices

$$S_W = 1/N_- \sum_{i \in \mathcal{Y}_{-1}} (\mathbf{x}_i - \mathbf{w}_{-1})(\mathbf{x}_i - \mathbf{w}_{-1})^\top$$

$$+ 1/N_+ \sum_{j \in \mathcal{Y}_{+1}} (\mathbf{x}_j - \mathbf{w}_{+1})(\mathbf{x}_j - \mathbf{w}_{+1})^\top$$

 Compute normal vector \mathbf{w}

$$\mathbf{w} = S_W^{-1}(\mathbf{w}_{+1} - \mathbf{w}_{-1})$$

 Compute threshold

$$\beta = 1/2 \mathbf{w}^\top (\mathbf{w}_{+1} + \mathbf{w}_{-1})$$

Output: \mathbf{w} , β

BCI with ML: Calibration and Feedback

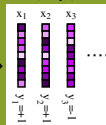
Calibration: continuous data

(markers provide information on mental states)



feature extraction

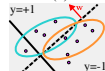
training data
(x_k, y_k)



classification
(training of the classifier)

optimizing parameters of the classifier f for: $f(x_k) \approx y_k$

(In LDA: $f(x) = w^T x + b$)



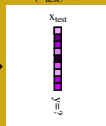
Feedback application: continuous data

(estimate mental state of most recent window)



feature extraction

'test' data
($x_{\text{test}}, ?$)



~~classification~~
(applying the classifier)

output
(prediction of the classifier)

$y = w^T x_{\text{test}} + b$



BCI Based on Event-Related Potentials (ERPs)

- User concentrates on a symbol
- Rows and columns are intensified randomly
- Target rows and columns elicit specific ERPs
- BCI detects target ERPs (averaged over few repetitions)

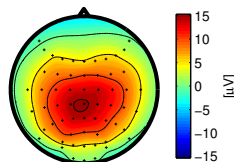
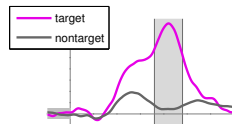
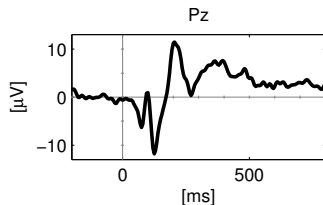
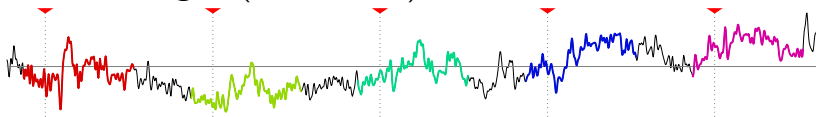


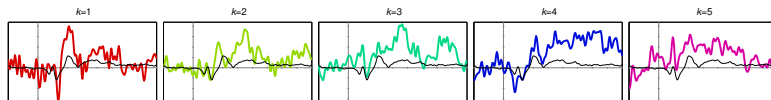
Illustration: Single-Trials and ERPs



Continuous Signal (with markers):

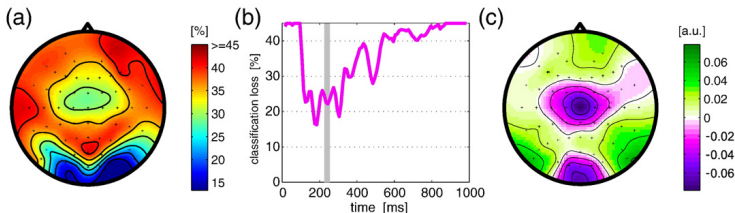


Segments (epochs) around stimulus markers:



35 / 40

Understanding the classifier



(a) Classification error on features from the time interval 115-535ms

(b) Classification error for intervals of 30ms duration

(c) Weight vector of classification on features from the time interval 220-250ms

[Blankertz et al., 2011]

Generalization and Model Evaluation

The goal of classification is **generalization**: Correct categorization/prediction of new data

How can we estimate generalization performance?

→ **Cross-validation**:

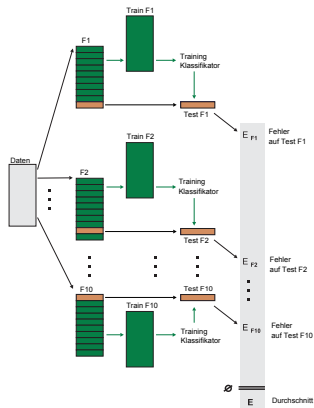
- Train model on part of data
- Test model on other part of data
- Repeat on different cross-validation *folds*
- Average performance on test set across all folds

Cross-Validation

Algorithm 1: Cross-Validation

Require: Data $(x_1, y_1) \dots, (x_N, y_N)$, Number of CV folds F

- 1: # Split data in F **disjunct** folds
- 2: **for** folds $f = 1, \dots, F$ **do**
- 3: # Train model on folds $\{1, \dots, F\} \setminus f$
- 4: # Compute prediction error on fold f
- 5: **end for**
- 6: # Average prediction error



Summary

Correlations between features can affect classification accuracy

Fisher proposed Linear Discriminant Analysis (LDA)

LDA maximizes *between class variance* while minimizing *within class variance*

If data is Gaussian with equal class covariances, then LDA is the optimal classifier

LDA is used in state-of-the-art BCI systems

We can use Cross-validation for Model Evaluation

References

- B. Blankertz, S. Lemm, M. Treder, S. Haufe, and K.-R. Müller. Single-trial analysis and classification of erp components—a tutorial. *Neuroimage*, 56(2):814–25, 2011. doi: 10.1016/j.neuroimage.2010.06.048.
- R. A. Fisher. The use of multiple measurements in taxonomic problems. *Annals of Eugenics*, 7:179–188, 1936.