

Pattern Recognition (PR)

Prof. Dr.-Ing. Andreas Maier
Pattern Recognition Lab (CS 5), Friedrich-Alexander-Universität Erlangen-Nürnberg
Winter Term 2020/21



This is a printable version of the slides of the lecture

Pattern Recognition (PR)
Winter term 2020/21
Friedrich-Alexander University of Erlangen-Nuremberg.

These slides are released under Creative Commons License Attribution CC BY 4.0.

Please feel free to reuse any of the figures and slides, as long as you keep a reference to the source of these slides at <https://lme.tf.fau.de/teaching/> acknowledging the authors Niemann, Hornegger, Hahn, Steidl, Nöth, Seitz, Rodriguez, Das and Maier.

Erlangen, January 8, 2021
Prof. Dr.-Ing. Andreas Maier

Discriminant Analysis II



Rank Reduced Linear Discriminant Analysis

Problem: How to choose an L -dimensional subspace with $L = K - 1$ that is good for LDA?

Idea: Maximize the spread of the L -dimensional projection of centroids.

Solution: Principal component analysis, i. e.
we compute the principal components of the
covariance matrix of the mean vectors

$$\mu'_y = \phi(\mu_y) \in \mathbb{R}^{K-1},$$

where $y = 1, 2, \dots, K$.

Rank Reduced Linear Discriminant Analysis (cont.)

In **Principle Component Analysis (PCA)** we compute a linear mapping $\Phi \in \mathbb{R}^{L \times (K-1)}$ that results in the highest spread of projected features:

$$\Phi^* = \operatorname{argmax}_{\Phi} \left(\frac{1}{K} \sum_{y=1}^K (\Phi \mu'_y - \Phi \bar{\mu}')^T (\Phi \mu'_y - \Phi \bar{\mu}') + \sum_{i=1}^L \lambda_i (\|\Phi_i\|_2^2 - 1) \right)$$

where we applied the **Lagrange multiplier** method to allow for the maximization of the spread subject to

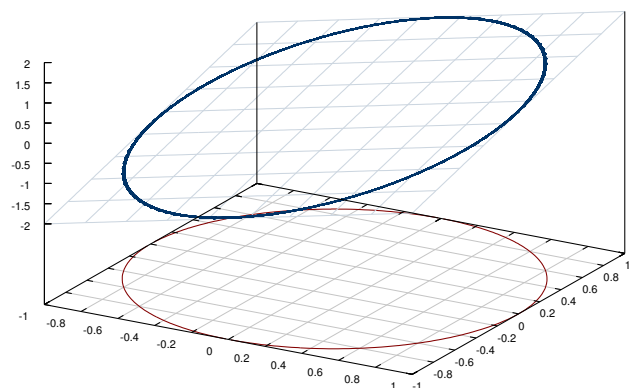
$$\|\Phi_i\|_2^2 = 1, \quad i = 1, \dots, K-1.$$

Here $\|\Phi_i\|_2^2$ denotes the L_2 norm of the i -th row vector of Φ .

Excursus: Optimization with Constraints

Lagrange Multipliers: simple example

- Find the maximum of
 $f(x, y) = x + y$
- constraint: $x^2 + y^2 = 1$
- Lagrange function:
 $L(x, y, \lambda) = x + y + \lambda(x^2 + y^2 - 1)$
- Set the partial derivatives to zero:
 $\frac{\partial L}{\partial x} = \frac{\partial L}{\partial y} = \frac{\partial L}{\partial \lambda} \stackrel{!}{=} 0$



Rank Reduced Linear Discriminant Analysis (cont.)

We need some facts from matrix calculus:


 www.matrixcookbook.com¹

1. Let μ denote the mean and Σ the covariance matrix of a random vector \mathbf{x} , then we get:

$$E[(\mathbf{A}\mathbf{x})^T(\mathbf{A}\mathbf{x})] = \text{tr}(\mathbf{A}\Sigma\mathbf{A}^T) + (\mathbf{A}\mu)^T(\mathbf{A}\mu)$$

2. The matrix derivative is:

$$\frac{\partial \text{tr}(\mathbf{X}\mathbf{B}\mathbf{X}^T)}{\partial \mathbf{X}} = \mathbf{X}\mathbf{B}^T + \mathbf{X}\mathbf{B}$$

¹ website currently off-line – you can still download it  here

Rank Reduced Linear Discriminant Analysis (cont.)

For our optimization problem this implies:

$$\begin{aligned} & \frac{\partial}{\partial \Phi} \left\{ \frac{1}{K} \sum_{y=1}^K (\Phi \mu'_y - \Phi \bar{\mu}')^T (\Phi \mu'_y - \Phi \bar{\mu}') + \sum_{i=1}^L \lambda_i (\|\Phi_i\|_2^2 - 1) \right\} \\ &= \frac{\partial}{\partial \Phi} \left\{ \frac{1}{K} \sum_{y=1}^K (\Phi (\mu'_y - \bar{\mu}'))^T (\Phi (\mu'_y - \bar{\mu}')) + \sum_{i=1}^L \lambda_i (\|\Phi_i\|_2^2 - 1) \right\} \\ &= \frac{\partial}{\partial \Phi} \left\{ \text{tr}(\Phi \Sigma_{\text{inter}} \Phi^T) + \sum_{i=1}^L \lambda_i (\|\Phi_i\|_2^2 - 1) \right\} \stackrel{!}{=} 0. \end{aligned}$$

Rank Reduced Linear Discriminant Analysis (cont.)

Now we compute the partial derivatives:

$$\frac{\partial}{\partial \Phi} \left\{ \text{tr}(\Phi \Sigma_{\text{inter}} \Phi^T) + \sum_{i=1}^L \lambda_i (\|\Phi_i\|_2^2 - 1) \right\} = 2\Phi \Sigma_{\text{inter}} + 2\lambda \Phi = 0$$

This results in the eigenvalue and eigenvector problem:

$$\Sigma_{\text{inter}} \Phi^T = \lambda' \Phi^T$$

Note:

In original PCA, the transform Φ maximizes the overall spread using the covariance matrix of all features:

$$\Sigma \Phi^T = \lambda' \Phi^T$$

Rank Reduced Linear Discriminant Analysis (cont.)

Input: training data: $S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), (\mathbf{x}_3, y_3), \dots, (\mathbf{x}_m, y_m)\}$

1. Compute the covariance matrix of transformed mean vectors

$$\hat{\Sigma}_{\text{inter}} = \frac{1}{K} \sum_{y=1}^K (\mu'_y - \bar{\mu}')(\mu'_y - \bar{\mu}')^T,$$

where $\bar{\mu}' = \frac{1}{K} \cdot \sum_{y=1}^K \mu'_y$.

2. Compute the L eigenvectors of the covariance matrix belonging to the largest eigenvalues.
3. The eigenvectors are the rows of the mapping Φ from the $(K - 1)$ - to the L -dimensional feature space.

Output: matrix Φ



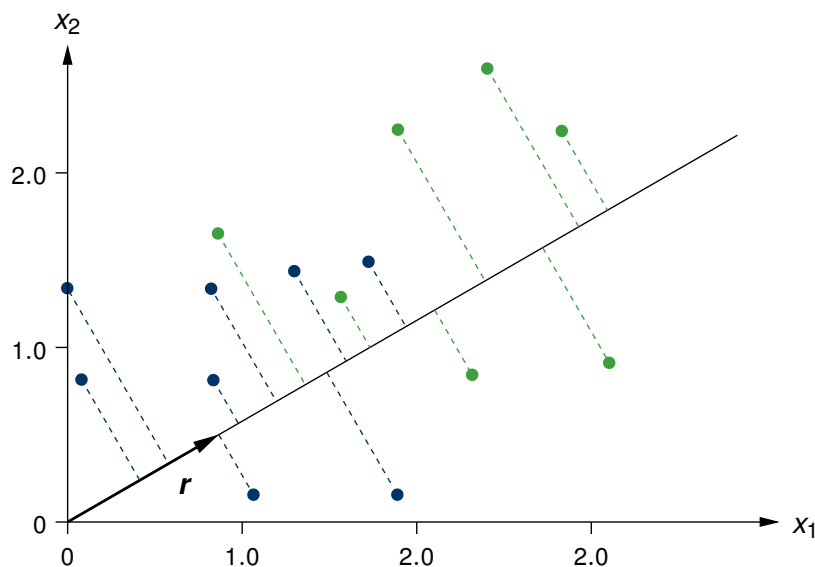
Next Time in Pattern Recognition



Fisher Transform

The described method to compute the LDA mapping is not the original derivation.

Original method



Fisher Transform

The described method to compute the LDA mapping is not the original derivation.

Original method

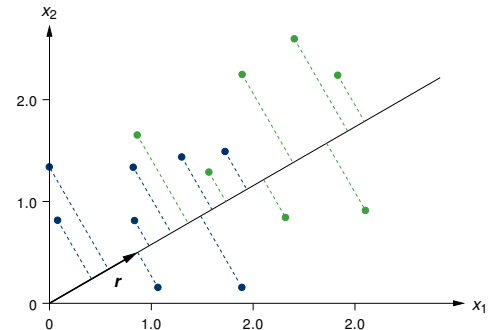
- Project samples \mathbf{x}_i onto a straight line with direction \mathbf{r} , $\|\mathbf{r}\|_2 = 1$:

$$\tilde{x}_i = \mathbf{x}_i^T \mathbf{r}$$

- Maximize the ratio of the between-class scatter and the within-class scatter:

$$\mathbf{r}^* = \underset{\mathbf{r}}{\operatorname{argmax}} J(\mathbf{r}) = \underset{\mathbf{r}}{\operatorname{argmax}} \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2}$$

- Classify by applying a threshold to \tilde{x}_i



Fisher Transform (cont.)

Finding \mathbf{r}^*

- Mean and scatter matrix for each class:

$$\boldsymbol{\mu}_k = \frac{1}{m_k} \sum_{\substack{i=1 \\ y_i=k}}^{m_k} \mathbf{x}_i$$

$$\mathbf{S}_k = \sum_{\substack{i=1 \\ y_i=k}}^{m_k} (\mathbf{x}_i - \boldsymbol{\mu}_k)(\mathbf{x}_i - \boldsymbol{\mu}_k)^T$$

- Within-class scatter matrix:

$$\mathbf{S}_W = \mathbf{S}_1 + \mathbf{S}_2$$

- Between-class scatter matrix:

$$\mathbf{S}_B = (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)(\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)^T$$

Fisher Transform (cont.)

Finding \mathbf{r}^*

4. Expressing $\tilde{\mu}_k$ and \tilde{s}_k^2 of the projected samples in terms of $\boldsymbol{\mu}_k$ and \mathbf{S}_k :

$$\begin{aligned}\tilde{\mu}_k &= \frac{1}{m_k} \sum_{\substack{i=1 \\ y_i=k}}^{m_k} \tilde{x}_i = \frac{1}{m_k} \sum_{\substack{i=1 \\ y_i=k}}^{m_k} \mathbf{r}^T \mathbf{x}_i = \mathbf{r}^T \boldsymbol{\mu}_k \\ \tilde{s}_k^2 &= \sum_{\substack{i=1 \\ y_i=k}}^{m_k} (\tilde{x}_i - \tilde{\mu}_k)^2 = \sum_{\substack{i=1 \\ y_i=k}}^{m_k} (\mathbf{r}^T \mathbf{x}_i - \mathbf{r}^T \boldsymbol{\mu}_k)^2 = \mathbf{r}^T \mathbf{S}_k \mathbf{r}\end{aligned}$$

5. Plug it into $J(\mathbf{r})$:

$$J(\mathbf{r}) = \frac{|\tilde{\mu}_1 - \tilde{\mu}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} = \frac{\mathbf{r}^T \mathbf{S}_B \mathbf{r}}{\mathbf{r}^T \mathbf{S}_W \mathbf{r}}$$

This is known as the **Generalized Rayleigh Quotient**.

Fisher Transform (cont.)

Finding \mathbf{r}^*

6. Maximizing the Generalized Rayleigh Quotient is equivalent to solving the following generalized eigenvalue problem:

$$\begin{aligned}\mathbf{S}_B \mathbf{r}^* &= \lambda \mathbf{S}_W \mathbf{r}^* \\ \mathbf{S}_W^{-1} \mathbf{S}_B \mathbf{r}^* &= \lambda \mathbf{r}^*\end{aligned}$$

7. **Note:** $\mathbf{S}_B \mathbf{r}^*$ is always in the direction of $\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2$;
no need to compute the eigenvalues and eigenvectors of $\mathbf{S}_W^{-1} \mathbf{S}_B$!

The direction of \mathbf{r}^* is:

$$\mathbf{r}^* = \mathbf{S}_W^{-1} (\boldsymbol{\mu}_1 - \boldsymbol{\mu}_2)$$

Fisher Transform (cont.)

- Usually the total linear mapping for LDA is computed dimension by dimension through the maximization of the Rayleigh ratio for each projection axis \mathbf{a}^* :

$$\mathbf{a}^* = \operatorname{argmax}_{\mathbf{a}} \frac{\mathbf{a}^T \Sigma_{\text{inter}} \mathbf{a}}{\mathbf{a}^T \Sigma_{\text{intra}} \mathbf{a}}$$

- The solution is a generalized eigenvalue problem: \mathbf{a} is the eigenvector of

$$\Sigma_{\text{intra}}^{-1} \Sigma_{\text{inter}}$$

that belongs to the largest eigenvalue.

Fisher Transform (cont.)

In literature the optimization problem is mostly rewritten:

- Equivalent constrained optimization problem

$$\begin{aligned} \text{maximize:} \quad & \mathbf{r}^T \Sigma_{\text{inter}} \mathbf{r} \\ \text{subject to:} \quad & \mathbf{r}^T \Sigma_{\text{intra}} \mathbf{r} = 1 \end{aligned}$$

- Lagrange multiplier method ($\lambda > 0$):

$$\mathbf{r}^* = \operatorname{argmax}_{\mathbf{r}} \{ \mathbf{r}^T \Sigma_{\text{inter}} \mathbf{r} - \lambda \mathbf{r}^T \Sigma_{\text{intra}} \mathbf{r} \}$$

Dimensionality Reduction

A few comments on dimensionality reduction:

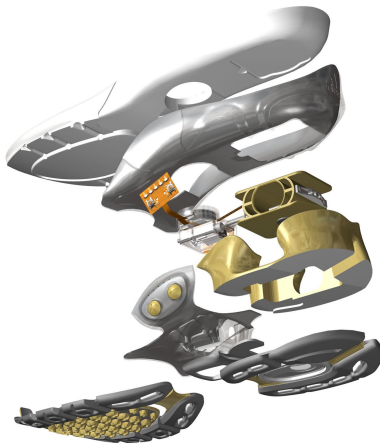
- PCA does not require a classified set of feature vectors (in contrast to LDA).
- PCA transformed features are approximately normally distributed (central limit theorem).
- Components of PCA transformed features are mutually independent.
- There exist many other methods for dimensionality reduction, e. g., Sammon transform, independent component analysis.
- Usually the estimation of transforms is computationally prohibited.
- **Johnson-Lindenstrauss lemma:** If vectors are projected onto a randomly selected subspace of suitably high dimension, then the distances between the vectors are approximately preserved.



Next Time in
Pattern Recognition



The adidas_1: A Digital Revolution in Sports



- For the first time ever, sport specific information can be processed with a running shoe
- A built-in microprocessor permits an adaptation of the shoe to the prevailing run situation
 - Running speed
 - Runner fatigue
 - Running surface
 - ...
- Pattern Recognition at the LME provides the algorithms used for recognition

The adidas_1: System Overview

Important parts of the adidas_1:

- A cushioning element (01) with a magnetic system for compression measurement
 - $f_{\text{sample}} = 1\text{kHz}$
 - resolution $\Delta d = 0.1\text{ mm}$
- A microcontroller and user interface (02)
 - $f_{\text{clock}} = 24\text{ MHz}$
 - 8 kB program memory
- A motor for cushioning adaptation using a cable system (03)



The adidas_1: Classification Framework Requirements

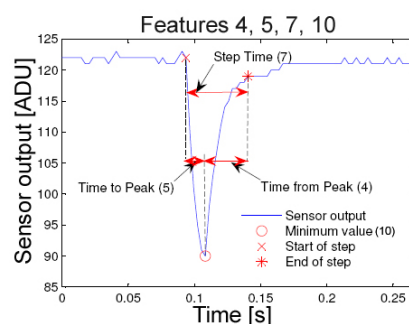
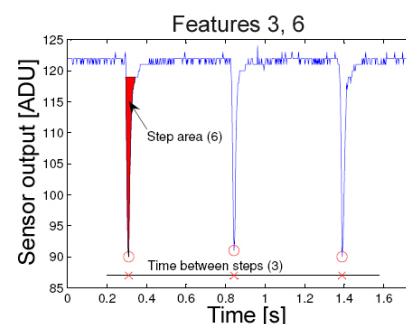
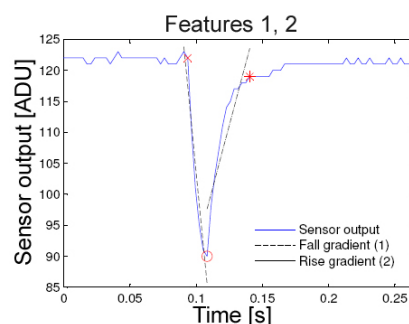
- Only a few, simple features can be calculated in real time
- The classification system has to be efficient, but computationally undemanding
- LDA classifier yields a linear decision boundary and can be implemented using a polynomial of order one with weights α_i and features x_i
- In the two class case:

$$\text{sgn}(\alpha^T \mathbf{x} + \alpha_0) = \text{sgn}(\alpha_1 x_1 + \alpha_2 x_2 + \dots + \alpha_d x_d + \alpha_0)$$

yields decision for either class.

Classification System: Computed Features

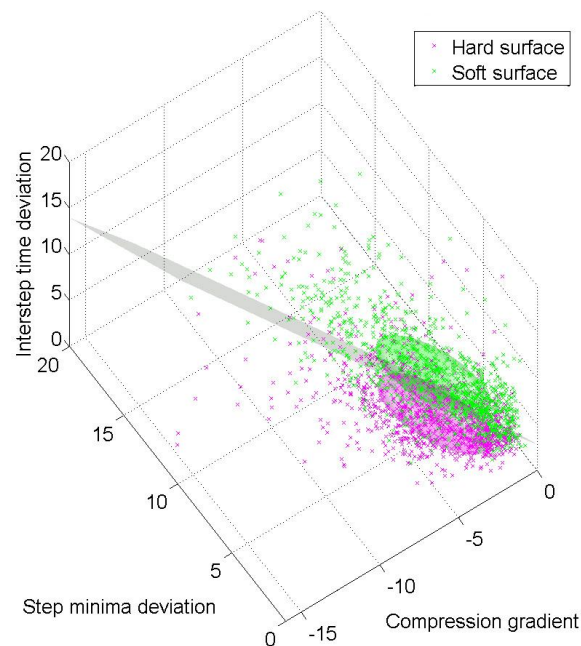
- 19 features initially computed for classification experiments
- Feature selection: 3 features selected for implementation



- Feature 8: mean value per step
- Feature 9: median value per step
- Feature 11: standard deviation (SD) of the values in one step
- Feature 12: SD of the minima
- Feature 13: SD of the means
- Feature 14: SD of the step standard deviation
- Feature 15: SD of the step time
- Feature 16: SD of the step area
- Feature 17: SD of the time between steps
- Feature 18: SD of the time to peak
- Feature 19: SD of the time from peak

Classification System: LDA Classifier Visualization

- Visualization of the decision region for hard/soft surface classification in 3D feature space



Shape Modeling

- Each shape is represented by n sampled surface points.
- Surface points are denoted by $\mathbf{p}_k \in \mathbb{R}^3$, $k = 1, 2, \dots, n$.
- The set of surface points is encoded in a single vector (shape vector):

$$\mathbf{x} = \begin{pmatrix} \mathbf{p}_1 \\ \mathbf{p}_2 \\ \vdots \\ \mathbf{p}_n \end{pmatrix} = \begin{pmatrix} p_{1,1} \\ p_{1,2} \\ p_{1,3} \\ p_{2,1} \\ \vdots \\ p_{n,3} \end{pmatrix} \in \mathbb{R}^{3n}$$

with $\mathbf{p}_k = (p_{k,1}, p_{k,2}, p_{k,3})^T$.

Shape Modeling (cont.)

We have m shapes, thus m shape vectors, and can generate the landmark configuration matrix:

$$\mathbf{L} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m]$$

Now we can compute the PCA of the columns of \mathbf{L} and get the spectral decomposition of the associated covariance matrix

$$\Sigma_L = \sum_i \lambda_i \mathbf{e}_i \mathbf{e}_i^T$$

where λ_i denote the eigenvalues and \mathbf{e}_i the eigenvectors.

Shape Modeling (cont.)

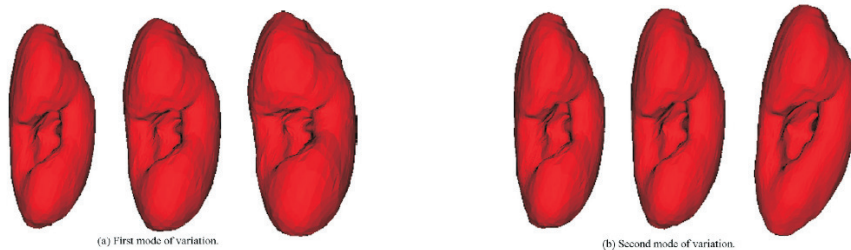
Shape vectors \mathbf{x}^* within the eigenvector space can be computed using linear combinations of I eigenvectors:

$$\mathbf{x}^* = \bar{\mathbf{x}} + \sum_{i=1}^I a_i \mathbf{e}_i$$

where $\bar{\mathbf{x}}$ denotes just the mean of the column vectors of \mathbf{L} and $a_i \in \mathbb{R}$ are the shape parameters.

Application of PCA: Segmentation

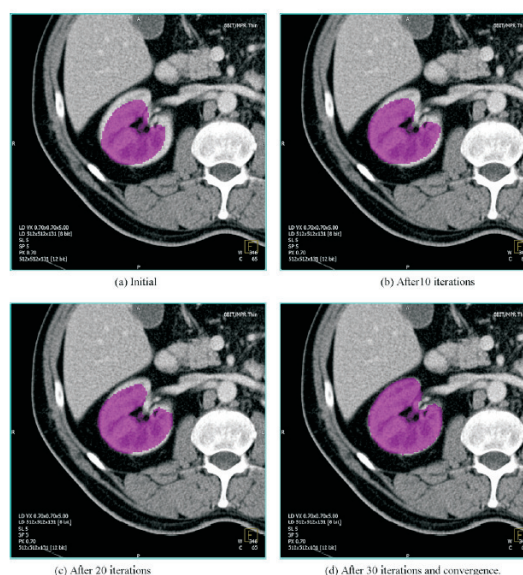
- Lung, liver or kidneys
- Generate and train an Active Shape Model (ASM) for such organs; requires training data and “gold standard” segmentation
- Once **point correspondences** are found, the different variations within the training data can be easily approximated by its Eigenvectors



M. Spiegel, D. Hahn, V. Daum, J. Wasza, J. Hornegger. "Segmentation of kidneys using a new active shape model generation technique based on non-rigid image registration", Computerized Medical Imaging and Graphics 2009

Fig.: Variation of the mean kidney shape along the first and second Eigenvector.

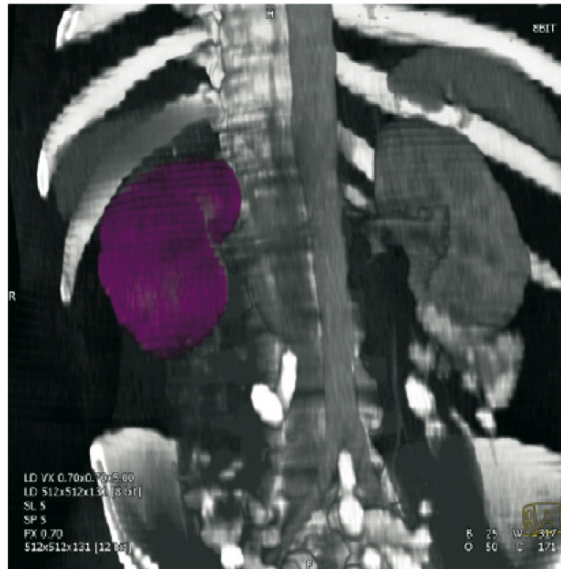
Application of PCA: Segmentation (cont.)



M. Spiegel, D. Hahn, V. Daum, J. Wasza, J. Hornegger. "Segmentation of kidneys using a new active shape model generation technique based on non-rigid image registration", Computerized Medical Imaging and Graphics 2009

Fig.: Iterative segmentation progress of a right kidney using an ASM.

Application of PCA: Segmentation (cont.)



M. Spiegel, D. Hahn, V. Daum, J. Wasza, J. Hornegger. "Segmentation of kidneys using a new active shape model generation technique based on non-rigid image registration", Computerized Medical Imaging and Graphics 2009

Fig.: 3-D view of the segmentation result.

Next Time in
Pattern Recognition



Notes on Regression

In the two class situation, we set $y \in \{-1, +1\}$ and use the decision rule:

$$y^* = \text{sgn}(\alpha^T \mathbf{x} + \alpha_0).$$

We can compute the linear decision boundary simply by least-square estimation.

Linear Regression (cont.)

For a given set of learning data we use matrix notation:

$$\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^T & 1 \\ \mathbf{x}_2^T & 1 \\ \vdots & \vdots \\ \mathbf{x}_m^T & 1 \end{pmatrix} \in \mathbb{R}^{m \times (d+1)} \quad \mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_m \end{pmatrix} \in \mathbb{R}^m$$

and define

$$\boldsymbol{\theta} = \begin{pmatrix} \boldsymbol{\alpha} \\ \alpha_0 \end{pmatrix}$$

Linear Regression (cont.)

One option to estimate θ is to solve the linear regression problem:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \|\mathbf{X}\theta - \mathbf{y}\|_2^2$$

Linear Regression (cont.)

The least-square estimator for the L_2 -norm:

$$\begin{aligned} \hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^m (\theta^T \mathbf{x}_i - y_i)^2 \\ &= \underset{\theta}{\operatorname{argmin}} (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y}) \end{aligned}$$

and thus we get

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

if the column vectors of \mathbf{X} are mutually independent.

Linear Regression (cont.)

A few obvious questions:

- Why should we prefer the Euclidean norm (L_2 -norm)?
- Will different norms lead to different results?
- Which norm and decision boundary is the best one?
- Can we incorporate prior knowledge in linear regression?

Ridge Regression

In *ridge regression* (also called *regularized regression*) we extend the objective function by an additional term constraining the Euclidean length of the parameter vector θ :

- It is linear regression with the log-likelihood penalized by $-\lambda \theta^T \theta$ where $\lambda > 0$, or alternatively
- It is extended by a prior distribution on the parameter vector θ

$$\theta = \mathcal{N}(0, \text{diag}(\tau^2))$$

Ridge Regression (cont.)

Regularized regression:

$$\begin{aligned}\hat{\theta} &= \underset{\theta}{\operatorname{argmin}} \|\mathbf{X}\theta - \mathbf{y}\|_2^2 + \lambda \|\theta\|_2^2 \\ &= \underset{\theta}{\operatorname{argmin}} (\mathbf{X}\theta - \mathbf{y})^T (\mathbf{X}\theta - \mathbf{y}) + \lambda \cdot \theta^T \theta\end{aligned}$$

and thus we get the estimator:

$$\hat{\theta} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$$

Ridge Regression (cont.)

Notes:

- The term $\lambda \mathbf{I}$ adds a positive constant λ to the diagonal elements.
- The problem is non-singular even if $\mathbf{X}^T \mathbf{X}$ is not of full rank.
- This was the main motivation of ridge regression when it was first introduced in statistics in 1970.
- The ridge solutions are not equivariant under scaling of the inputs:
standardize the input before solving the regression problem!
- The intercept α_0 should not be penalized:
 - Center the input \mathbf{x}_i .
 - Estimate α_0 by $\bar{y} = \frac{1}{m} \sum_{i=1}^m y_i$.
 - Estimate the remaining coefficients by a ridge regression without intercept.
Matrix \mathbf{X} has d columns (instead of $d + 1$).

Ridge Regression (cont.)

Statistical approach: parameters α_j are random variables

- Suppose

$$\forall 1 \leq i \leq m: \quad y_i \sim \mathcal{N}(\underbrace{\alpha^T \mathbf{x}_i + \alpha_0}_{\text{mean}}, \underbrace{\sigma^2}_{\text{variance}}) = \frac{1}{\sqrt{2\pi} \cdot \sigma} e^{-\frac{1}{2} \cdot \frac{(y_i - \alpha^T \mathbf{x}_i - \alpha_0)^2}{\sigma^2}}$$

- Parameters α_j are assumed to be independent of each other.
- Prior distribution of α_j :

$$\forall 1 \leq j \leq d: \quad \alpha_j \sim \mathcal{N}(0, \tau^2) = \frac{1}{\sqrt{2\pi} \cdot \tau} e^{-\frac{1}{2} \cdot \frac{(\alpha_j - 0)^2}{\tau^2}}$$

Ridge Regression (cont.)

- Maximizing the posterior probability of α for given σ^2 and τ^2 :

$$\begin{aligned} \operatorname{argmax}_{\alpha} \prod_{i=1}^m p(\alpha | y_i) &= \operatorname{argmax}_{\alpha} \left\{ \prod_{i=1}^m p(\alpha) \cdot p(y_i | \alpha) \right\} \\ &= \operatorname{argmax}_{\alpha} \left\{ \prod_{j=1}^d p(\alpha_j) \cdot \prod_{i=1}^m p(y_i | \alpha) \right\} \\ &= \operatorname{argmax}_{\alpha} \left\{ \sum_{j=1}^d \log p(\alpha_j) + \sum_{i=1}^m \log p(y_i | \alpha) \right\} \\ &= \operatorname{argmax}_{\alpha} \left\{ -\frac{1}{2\tau^2} \sum_{j=1}^d \alpha_j^2 - \frac{1}{2\sigma^2} \sum_{i=1}^m (y_i - \alpha^T \mathbf{x}_i - \alpha_0)^2 \right\} \\ &= \operatorname{argmin}_{\alpha} \left\{ \frac{\sigma^2}{\tau^2} \sum_{j=1}^d \alpha_j^2 + \sum_{i=1}^m (y_i - \alpha^T \mathbf{x}_i - \alpha_0)^2 \right\} \end{aligned}$$

Ridge Regression (cont.)

- Maximizing the posterior probability of α for given σ^2 and τ^2 :

$$\operatorname{argmax}_{\alpha} \prod_{i=1}^m p(\alpha | y_i) = \operatorname{argmin}_{\alpha} \left\{ \lambda \alpha^T \alpha + (\mathbf{X}\alpha - \mathbf{y})^T (\mathbf{X}\alpha - \mathbf{y}) \right\} \quad \text{with} \quad \lambda = \frac{\sigma^2}{\tau^2}$$

- The ridge estimate is the mode of the posterior pdf!

Lasso

Regularized regression using a mixture of L_2 - and L_1 -norm, where the residual is penalized using the L_2 -norm and the regularizer uses the L_1 -norm:

$$\hat{\boldsymbol{\theta}} = \operatorname{argmin}_{\boldsymbol{\theta}} \|\mathbf{X}\boldsymbol{\theta} - \mathbf{y}\|_2^2 + \lambda \cdot \|\boldsymbol{\theta}\|_1$$

The lasso is used to compute a sparse solution of the system of linear equations, i. e. the number of non-zero elements in $\boldsymbol{\theta}$ shall be small.

Lessons Learned

- Principal component analysis
- Linear discriminant analysis with and without dimension reduction
- Both PCA and LDA relate to an eigenvalue eigenvector problem
- Alternative formulation of LDA using the Fisher transform
- Linear and ridge regression for classification



Next Time in Pattern Recognition



Further Readings

You are required to be familiar with **linear algebra** and **matrix calculus**:

- SIAMS best selling book in the last decade:
Lloyd N. Trefethen, David Bau III:
Numerical Linear Algebra,
SIAM, Philadelphia, 1997.
- All about matrix derivatives and related problems is described in the Matrix Cookbook: <http://www.matrixcookbook.com>

Basics on **discriminant analysis** can be found in

- T. Hastie, R. Tibshirani, and J. Friedman:
**The Elements of Statistical Learning –
Data Mining, Inference, and Prediction**,
2nd edition, Springer, New York, 2009.

Further Readings (cont.)

Details on the adidas_1 shoe and the implemented classifier:

- B. Eskofier, F. Hönig, P. Kühner:
Classification of Perceived Running Fatigue in Digital Sports,
Proceedings of the 19th International Conference on Pattern Recognition (ICPR
2008), Tampa, Florida, U. S. A., 2008

Details on the shape modeling of kidneys and its application to segmentation:

- M. Spiegel, D. Hahn, V. Daum, J. Wasza, J. Hornegger:
**Segmentation of kidneys using a new active shape model generation technique
based on non-rigid image registration**,
Computerized Medical Imaging and Graphics 2009 33(1):29-39

Comprehensive Questions

- What is the difference between PCA and LDA?
- How can PCA and LDA be combined to achieve a high rank reduction?
- Write down a straight forward objective function for linear regression!
- What happens if we replace the L_2 -norm by another norm?