

Cognitive Algorithms - Exercise 2

Daniel Wujecki

May 9, 2020

- 1 Task 1 - Linear Classification with NCC and Perceptron
- 2 Task 2 - Covariance Matrices
- 3 Task 3 - Linear Discriminant Analysis
- 4 Task 4 - Understanding LDA: the whitening operation
- 5 Model evaluation

1 Task 1 - Linear Classification with NCC and Perceptron

2 Task 2 - Covariance Matrices

3 Task 3 - Linear Discriminant Analysis

4 Task 4 - Understanding LDA: the whitening operation

5 Model evaluation

Classification

- Common problem in machine learning
- find a mapping $f : \mathbb{R}^d \rightarrow \mathcal{C}$ that estimates the class $y \in \mathcal{C}$ for a set of features $\mathbf{x} \in \mathbb{R}^d$

Classification

- Common problem in machine learning
- find a mapping $f : \mathbb{R}^d \rightarrow \mathcal{C}$ that estimates the class $y \in \mathcal{C}$ for a set of features $\mathbf{x} \in \mathbb{R}^d$
- binary classification with $\mathcal{C} = \{-1, 1\}$
- simple mapping $f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - \beta)$

Classification

- Common problem in machine learning
- find a mapping $f : \mathbb{R}^d \rightarrow \mathcal{C}$ that estimates the class $y \in \mathcal{C}$ for a set of features $\mathbf{x} \in \mathbb{R}^d$
- binary classification with $\mathcal{C} = \{-1, 1\}$
- simple mapping $f(\mathbf{x}) = \text{sign}(\mathbf{w}^T \mathbf{x} - \beta)$
- $\mathbf{w} \in \mathbb{R}^d$ is called the weight vector
- β is the bias or offset

Supervised learning

- How do we find a mapping f that estimates the *correct* labels $y \in \mathcal{C}$?

Supervised learning

- How do we find a mapping f that estimates the *correct* labels $y \in \mathcal{C}$?
- usually we need a set of observations with known class labels
- training data set can be given as a matrix $X \in \mathbb{R}^{d \times n}$ and a vector $y \in \mathcal{C}^n$

$$X = \begin{bmatrix} 0.1 & -3.2 & \cdots & 2.1 & 1.4 \\ \vdots & & & \vdots & \\ 0.1 & -3.2 & \cdots & 2.1 & 1.4 \end{bmatrix}$$
$$\mathbf{y} = [-1 \ 1 \ \cdots \ 1 \ -1]$$

Supervised learning

- How do we find a mapping f that estimates the *correct* labels $y \in \mathcal{C}$?
- usually we need a set of observations with known class labels
- training data set can be given as a matrix $X \in \mathbb{R}^{d \times n}$ and a vector $y \in \mathcal{C}^n$

$$X = \begin{bmatrix} 0.1 & -3.2 & \cdots & 2.1 & 1.4 \\ \vdots & & & \vdots & \\ 0.1 & -3.2 & \cdots & 2.1 & 1.4 \end{bmatrix}$$
$$\mathbf{y} = [-1 \ 1 \ \cdots \ 1 \ -1]$$

- choose f such that the number of misclassified observations is minimized

Repetition Perceptron

- algorithm to reduce the classification error iterative with gradient descent
- include the bias β into w (see first lecture) $\tilde{w} = \begin{pmatrix} \beta \\ w_1 \\ \vdots \\ w_d \end{pmatrix}$ $\tilde{x} = \begin{pmatrix} 1 \\ x_1 \\ \vdots \\ x_d \end{pmatrix}$ $w^T x - \beta = \tilde{w}^T \tilde{x}$

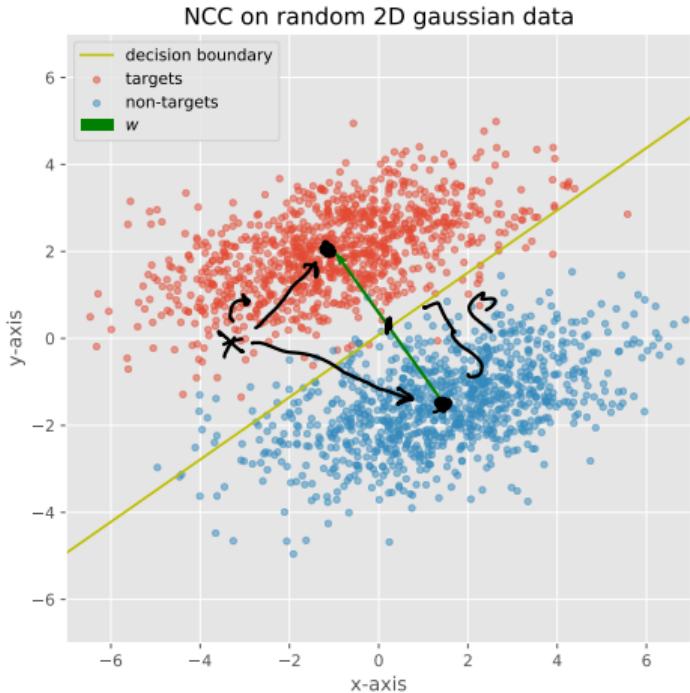
Repetition Perceptron

- algorithm to reduce the classification error iterative with gradient descent
- include the bias β into \mathbf{w} (see first lecture)

- 1 initialize \mathbf{w} randomly
- 2 until convergence (runtime constraint or no misclassifications)
 - a pick a random misclassified point \mathbf{x}_k
 - b update the weights with gradient descent: $\mathbf{w}_{new} = \mathbf{w}_{old} - \eta \mathbf{x}_k y_k$



Repetition Nearest Centroid Classifier



- compare distance of data points x to class means
- decide for class with closer center

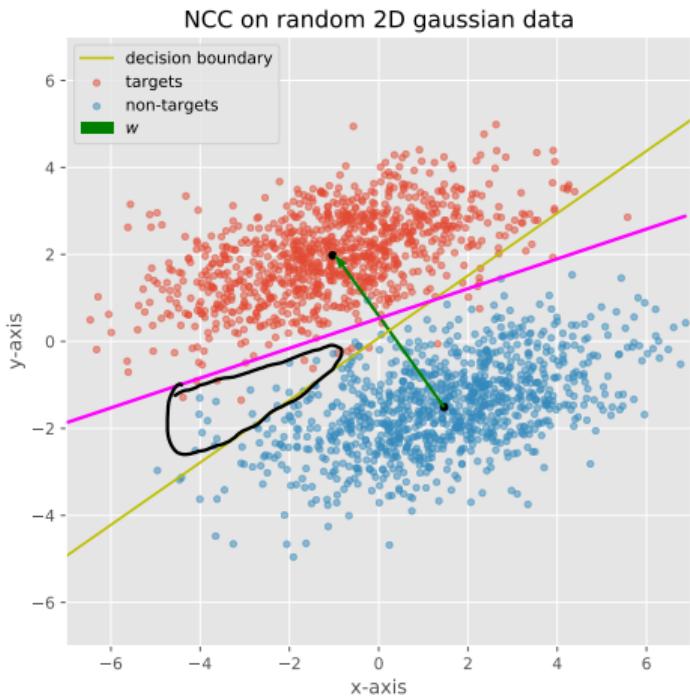
$$\mathbf{w} = \mathbf{w}_o - \mathbf{w}_\Delta$$

$$\beta = \frac{1}{2}(\mathbf{w}_o^T \mathbf{w}_o - \mathbf{w}_\Delta^T \mathbf{w}_\Delta)$$

$$\rightarrow f(x) = \text{sign}(\mathbf{w}^T \mathbf{x} - \beta) \leftarrow$$

$$\boxed{1} \dots \boxed{1} \approx \boxed{1}$$

Repetition Nearest Centroid Classifier



- compare distance of data points x to class means
- decide for class with closer center

$$\mathbf{w} = \mathbf{w}_o - \mathbf{w}_\Delta$$

$$\beta = \frac{1}{2}(\mathbf{w}_o^T \mathbf{w}_o - \mathbf{w}_\Delta^T \mathbf{w}_\Delta)$$

$$f(x) = sign(\mathbf{w}^T x - \beta)$$

- correlation not considered, thus often not optimal

Task 1: Bias term of the NCC

Show that

$$\beta = \frac{1}{2} (\underbrace{\mathbf{w}_o^T \mathbf{w}_o - \mathbf{w}_\Delta^T \mathbf{w}_\Delta}_{\text{underbrace}}) = \underbrace{\mathbf{w}^T \left(\frac{\mathbf{w}_o + \mathbf{w}_\Delta}{2} \right)}_{\text{underbrace}}$$

$$\mathbf{w}^T \left(\frac{\mathbf{w}_o + \mathbf{w}_\Delta}{2} \right) = \mathbf{w}^T \left(\frac{1}{2} (\mathbf{w}_o + \mathbf{w}_\Delta) \right) = \frac{1}{2} \mathbf{w}^T (\mathbf{w}_o + \mathbf{w}_\Delta)$$

$$= \frac{1}{2} (\mathbf{w}_o - \mathbf{w}_\Delta)^T (\mathbf{w}_o + \mathbf{w}_\Delta)$$

$$= \frac{1}{2} \left(\mathbf{w}_o^T \mathbf{w}_o + \underbrace{\mathbf{w}_o^T \mathbf{w}_\Delta}_{\cancel{\mathbf{w}_o^T \mathbf{w}_\Delta}} - \mathbf{w}_\Delta^T \mathbf{w}_o - \mathbf{w}_\Delta^T \mathbf{w}_\Delta \right)$$

$$= \frac{1}{2} (\mathbf{w}_o^T \mathbf{w}_o - \mathbf{w}_\Delta^T \mathbf{w}_\Delta)$$

$$\begin{cases} (\mathbf{v})^T \mathbf{w} = \underline{\mathbf{v}^T \mathbf{w}} \\ \mathbf{v}^T \mathbf{w} = \mathbf{w}^T \mathbf{v} \end{cases}$$

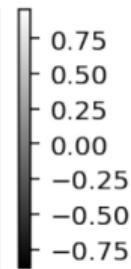
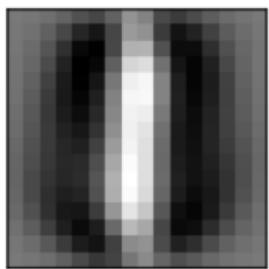
Task 1: Bias term of the NCC

Assignment 1: High accuracy but "lousy" weight vector?

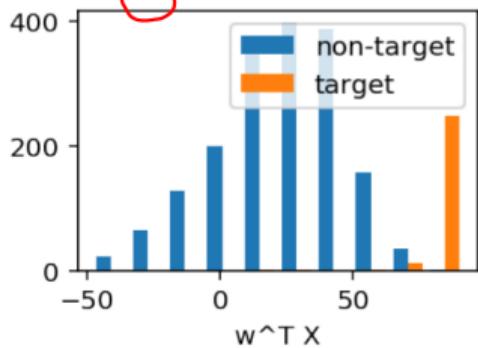
NCC

$$\omega = \omega_0 - \omega_\Delta$$

NCC



Acc 93.32336821126059%

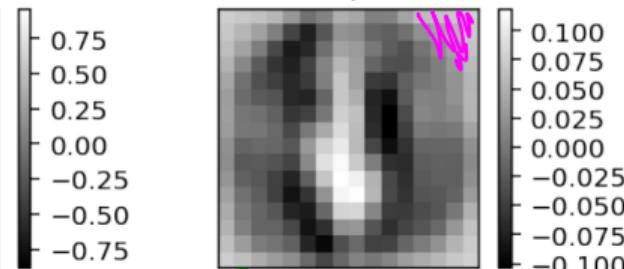


Perceptron

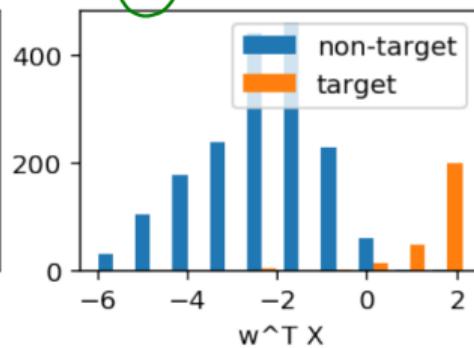
Perceptron

$$\omega$$

$$\vdots$$



Acc 98.50523168908819%



Assignment 1: High accuracy but "lousy" weight vector?

- purpose of weight vector is two-fold:
 - amplify important features
 - suppress unimportant features
- Example with some noise ϵ

Assignment 1: High accuracy but "lousy" weight vector?

- purpose of weight vector is two-fold:
 - amplify important features
 - suppress unimportant features
- Example with some noise ϵ

$$\begin{aligned} \omega^T x &= \begin{pmatrix} y + \epsilon \\ \epsilon \end{pmatrix}^\top \begin{pmatrix} 1 \\ -1 \end{pmatrix} \\ x &= \begin{bmatrix} y + \epsilon \\ \epsilon \end{bmatrix} = y + \epsilon - \epsilon = y \\ w &= \begin{bmatrix} 1 \\ -1 \end{bmatrix} \\ w^T x - \beta &= y + \epsilon - \epsilon - 0 = y \end{aligned}$$

- correct classification of class label $y \in \mathcal{C}$
 → w is the optimal weight vector for this data sample

- 1 Task 1 - Linear Classification with NCC and Perceptron
- 2 Task 2 - Covariance Matrices
- 3 Task 3 - Linear Discriminant Analysis
- 4 Task 4 - Understanding LDA: the whitening operation
- 5 Model evaluation

Covariance Matrix

$$x_1, \dots, x_n \in \mathbb{R}^d$$

- let X_1, \dots, X_d be random variables (can hold the different values of d features)
- their covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ is defined such that

$$\Sigma_{i,j} = \Sigma_{j,i} = \mathbb{E} [(X_i - \mathbb{E}(X_i))(X_j - \mathbb{E}(X_j))]$$

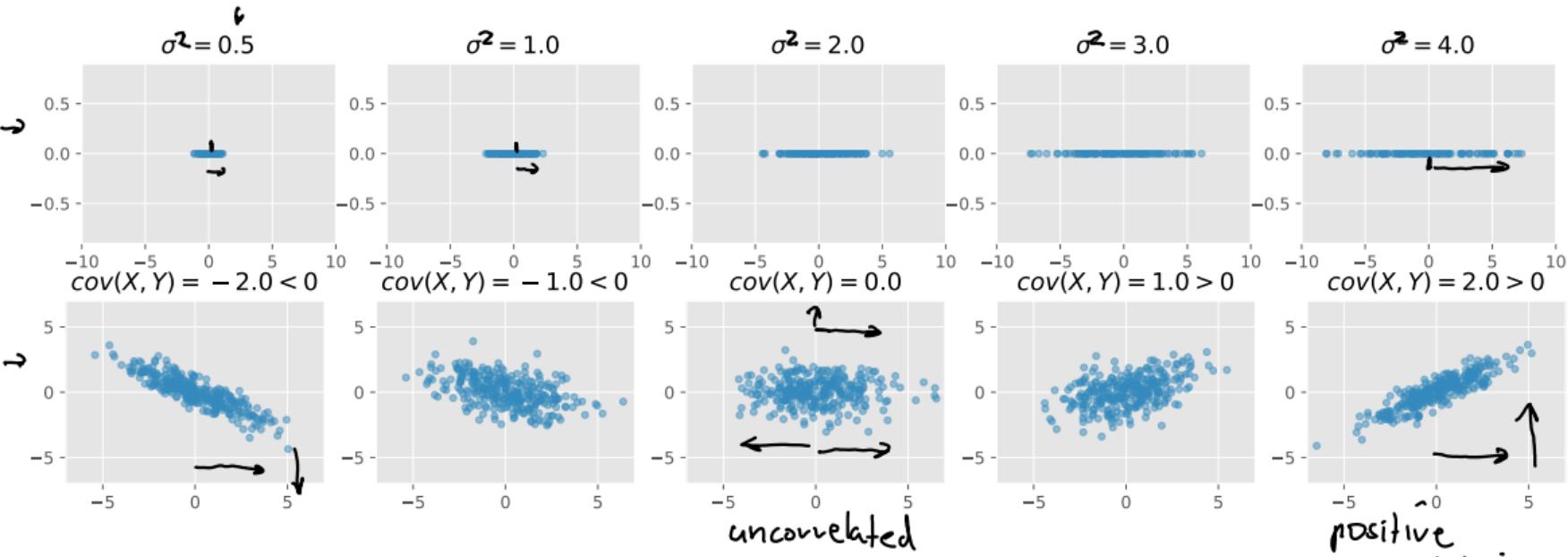
Covariance Matrix

- let X_1, \dots, X_d be random variables (can hold the different values of d features)
- their covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ is defined such that

$$\Sigma_{i,j} = \Sigma_{j,i} = \mathbb{E} [(X_i - \mathbb{E}(X_i))(X_j - \mathbb{E}(X_j))]$$

- covariance measures linear relationship between two random variables
- for $i = j$ (the diagonal) we get $\text{var}(X_i)$ (expected squared deviation from mean)

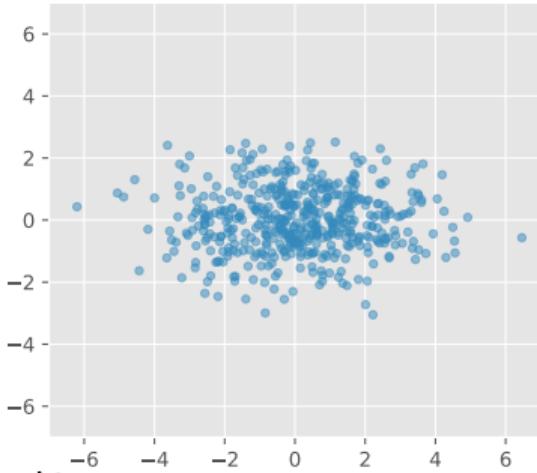
Covariance Matrix



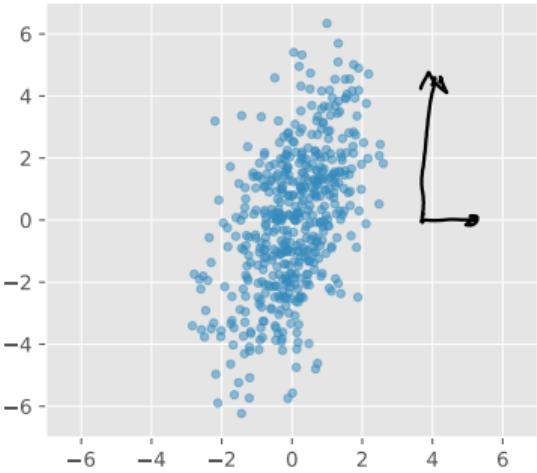
- covariance measures linear relationship between two random variables
- for $i = j$ (the diagonal) we get $\text{var}(X_i)$ (expected squared deviation from mean)

Covariance Matrix - Quiz

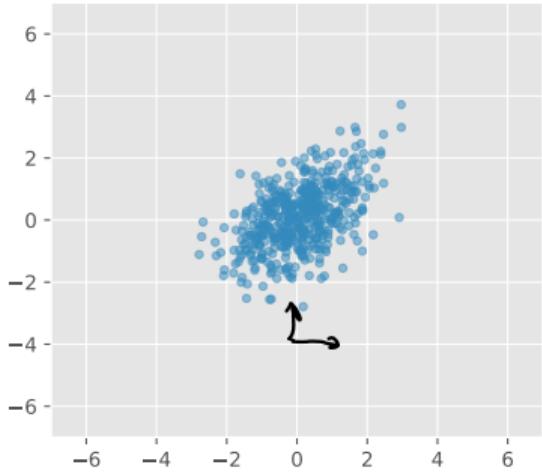
$$\Sigma_? \quad \xi_3$$



$$\Sigma_? \quad \xi_1$$



$$\xi_4$$



$$\xi_1^k = \xi_1$$

$$\Sigma_1 = \begin{pmatrix} 1 & 1 \\ 1 & 5 \end{pmatrix},$$

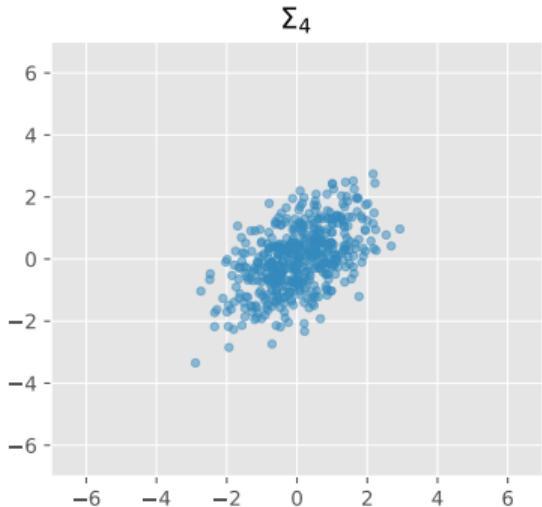
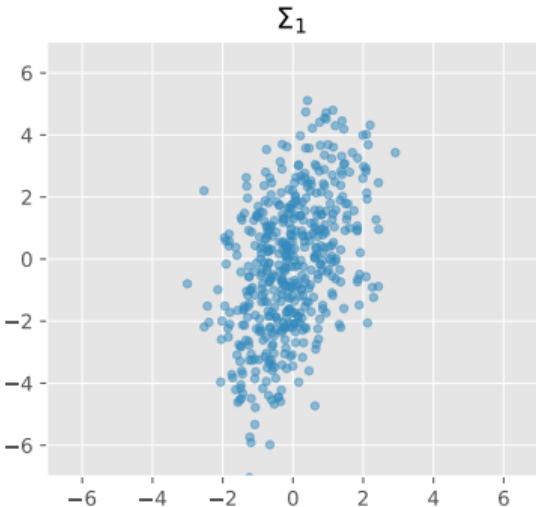
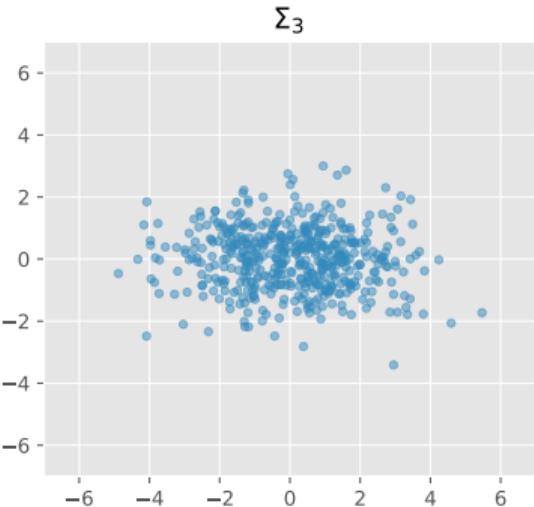
$$\Sigma_2 = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix},$$

$$\Sigma_3 = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix},$$

$$\Sigma_4 = \begin{pmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix},$$

$$\Sigma_5 = \begin{pmatrix} 5 & 1 \\ 1 & 1 \end{pmatrix}$$

Covariance Matrix - Quiz Solution



$$\Sigma_1 = \begin{pmatrix} 1 & 1 \\ 1 & 5 \end{pmatrix}, \quad \Sigma_2 = \begin{pmatrix} 2 & 2 \\ 2 & 2 \end{pmatrix}, \quad \Sigma_3 = \begin{pmatrix} 3 & 0 \\ 0 & 1 \end{pmatrix}, \quad \Sigma_4 = \begin{pmatrix} 1 & \frac{1}{2} \\ \frac{1}{2} & 1 \end{pmatrix}, \quad \Sigma_5 = \begin{pmatrix} 5 & 1 \\ 1 & 1 \end{pmatrix}$$

Empirical Covariance Matrix

- let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ be n observations of d random variables
- let be $\bar{\mathbf{x}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$ the empirical mean
- the empirical covariance matrix is defined as

$$\Sigma_X = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T$$

$$\begin{pmatrix} \vdots \\ \vdots \\ \vdots \end{pmatrix} \begin{bmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{bmatrix} \simeq \begin{pmatrix} \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot \end{pmatrix}$$

Empirical Covariance Matrix

- let $\mathbf{x}_1, \dots, \mathbf{x}_n \in \mathbb{R}^d$ be n observations of d random variables
- let be $\bar{\mathbf{x}} = \frac{1}{n} \sum_{k=1}^n \mathbf{x}_k$ the empirical mean
- the empirical covariance matrix is defined as

$$\Sigma_X = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T$$

- let $X = (\mathbf{x}_1 - \bar{\mathbf{x}}, \dots, \mathbf{x}_n - \bar{\mathbf{x}}) \in \mathbb{R}^{d \times n}$ denote the matrix of the centered data, then

$$\Sigma_X = \frac{1}{n} X X^T$$

Task 2.1

Compute the empirical covariance matrix for the following four data points

$$X = \begin{bmatrix} -1 & -1 & 1 & 1 \\ -1 & 0 & 0 & 1 \end{bmatrix} \in \mathbb{R}^{2,4}$$

$$\bar{x} = \frac{1}{4} \begin{bmatrix} -1 - 1 + 1 + 1 \\ -1 + 0 + 0 + 1 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 0 \\ 0 \end{bmatrix}$$

$$\Sigma_X = \frac{1}{n} X X^T = \frac{1}{4} \begin{bmatrix} -1 & -1 & 1 & 1 \\ -1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} -1 & -1 \\ -1 & 0 \\ 1 & 0 \\ 1 & 1 \end{bmatrix} = \frac{1}{4} \begin{bmatrix} 4 & 2 \\ 2 & 2 \end{bmatrix}$$

Task 2.1 - practice with numpy

→ see jupyter notebook

Task 2.1 - practice with numpy

→ see jupyter notebook

- unbiased estimation of the covariance matrix

$$\Sigma_X = \frac{1}{n-1} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T$$

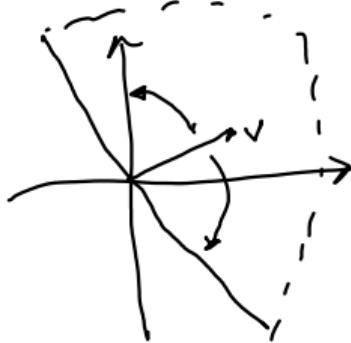
- instead of

$$\Sigma_X = \frac{1}{n} \sum_{k=1}^n (\mathbf{x}_k - \bar{\mathbf{x}})(\mathbf{x}_k - \bar{\mathbf{x}})^T$$

- for big n better estimation of the covariance
- based on statistics methods

Task 2.2 - properties of covariance matrices

- covariance matrices are always symmetric $\Sigma_X^T = \Sigma_X$
- covariance matrices are (semi-)positive definite, that is, that for every vector $v \in \mathbb{R}^d \setminus \{0\}$ the following inequality holds



$$\begin{aligned} v^T \Sigma_X v &\geq 0 \\ v^T v &\geq 0 \end{aligned}$$

$$\begin{aligned} v^T v &= \underbrace{\cos(\alpha)}_{\geq 0} \underbrace{\|v\| \|v\|}_{\geq 0} \\ \alpha &\in [-\frac{\pi}{2}, \frac{\pi}{2}] \end{aligned}$$

Task 2.2 - properties of covariance matrices

- covariance matrices are always symmetric $\Sigma_X^T = \Sigma_X$
- covariance matrices are (semi-)positive definite, that is, that for every vector $\mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\}$ the following inequality holds

$$\mathbf{v}^T \Sigma_X \mathbf{v} \geq 0$$

- this is equivalent to the observation $\lambda \geq 0$ for all eigenvalues λ of Σ_X

Task 2.2 - properties of covariance matrices

Proof I: all eigenvalues $\lambda \geq 0 \Rightarrow \forall v \in \mathbb{R}^d \setminus \{0\} : v^T \Sigma_X v \geq 0$

$$\Sigma_X = U \Lambda U^T$$

$$x \in \mathbb{R}^d \setminus \{0\} : x^T \Sigma_X x = x^T U \Lambda U^T x \\ = y^T y$$

$$\begin{array}{c} y = U^T x \\ y^T = (U^T x)^T = x^T U \\ \hline \boxed{Uy = x} \end{array}$$

$$= [y_1 \dots y_d] \begin{bmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & & \lambda_d \end{bmatrix} \begin{bmatrix} x_1 \\ \vdots \\ x_d \end{bmatrix}$$

$$= [y_1 \dots y_d] \begin{pmatrix} y_1 \lambda_1 \\ \vdots \\ y_d \lambda_d \end{pmatrix}$$

$$x^T \Sigma_X x \geq 0 \quad \Leftrightarrow \quad y^T y = y_1^2 \underline{\lambda_1} + y_2^2 \underline{\lambda_2} + \dots + y_d^2 \underline{\lambda_d} \geq 0$$

Task 2.2 - properties of covariance matrices

Proof II: $\forall \mathbf{v} \in \mathbb{R}^d \setminus \{\mathbf{0}\} : \mathbf{v}^T \Sigma_X \mathbf{v} \geq 0 \Rightarrow$ all eigenvalues $\lambda \geq 0$

$$\Sigma_X \mathbf{u} = \mathbf{u} \lambda \quad \underbrace{\mathbf{u}^T \Sigma_X \mathbf{u}}_{\geq 0} = \mathbf{u}^T (\lambda \mathbf{u}) = \lambda \mathbf{u}^T \mathbf{u} = \lambda \underbrace{\|\mathbf{u}\|_2^2}_{\geq 0}$$

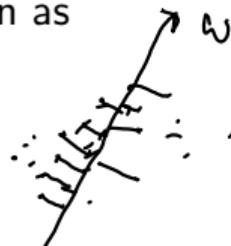
$$\mathbf{v}^T \Sigma_X \mathbf{v} > 0 ; \lambda_i > 0 \quad \text{positive definite}$$

$$\mathbf{v}^T \Sigma_X \mathbf{v} \geq 0 ; \lambda_i \geq 0 \quad (\text{semi}) \text{PD}$$

Task 2.3 - Variance in a certain direction

The variance of centered data $x_1, \dots, x_n \in \mathbb{R}^d$ in direction $w \in \mathbb{R}^d$ is given as

$$\frac{1}{n} \sum_{k=1}^n (w^T x_k)^2.$$



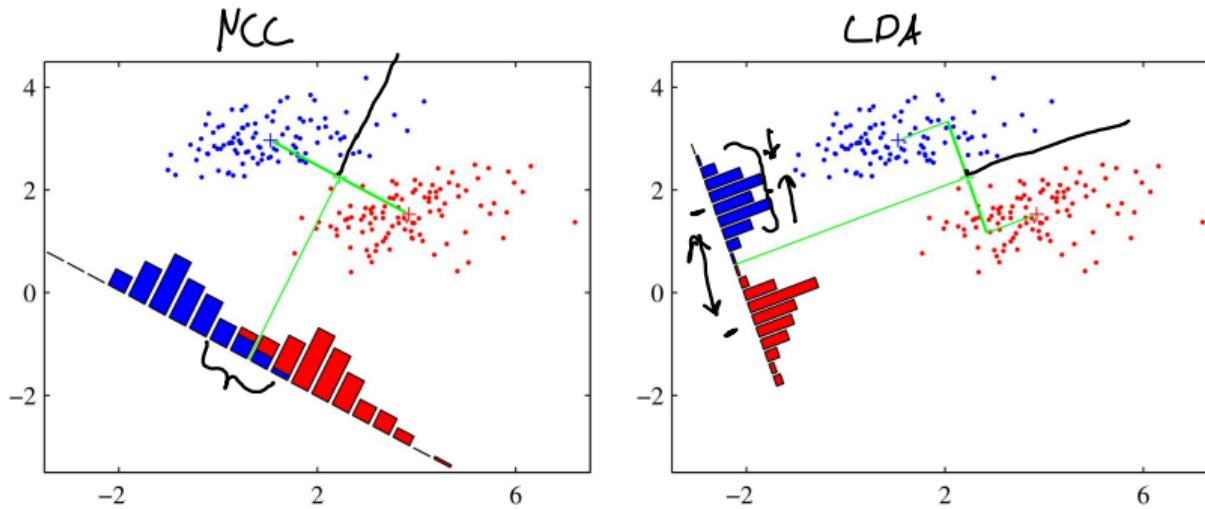
We show that this is equal to

$$\begin{aligned} \frac{1}{n} \sum_{k=1}^n (w^T x_k)^2 &= \frac{1}{n} \sum_{k=1}^n (w^T x_k)(x_k^T w) = w^T \underbrace{\left(\frac{1}{n} \sum_{k=1}^n x_k x_k^T \right)}_{\Sigma_X} w \\ &= w^T \Sigma_X w \end{aligned}$$

Task 2.3 - Variance in a certain direction

- 1 Task 1 - Linear Classification with NCC and Perceptron
- 2 Task 2 - Covariance Matrices
- 3 Task 3 - Linear Discriminant Analysis
- 4 Task 4 - Understanding LDA: the whitening operation
- 5 Model evaluation

Linear Discriminant Analysis



LDA gives us a weight vector w that

- maximizes the mean class difference
- minimizes the variance in each class

Linear Discriminant Analysis

- given class means $\mathbf{w}_o, \mathbf{w}_\Delta \in \mathbb{R}^d$ and number of points n_o, n_Δ per class

$$\mathbf{w} = \Sigma_X^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta)$$

$$\beta = \underbrace{\mathbf{w}^T \left(\frac{\mathbf{w}_o + \mathbf{w}_\Delta}{2} \right)}_{\log(1) = 0} \left[+ \log \left(\frac{n_o}{n_\Delta} \right) \right]$$

$$n_o = n_\Delta$$

Linear Discriminant Analysis

- given class means $\mathbf{w}_o, \mathbf{w}_\Delta \in \mathbb{R}^d$ and number of points n_o, n_Δ per class

$$\mathbf{w} = \Sigma_X^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta)$$

$$\beta = \mathbf{w}^T \left(\frac{\mathbf{w}_o + \mathbf{w}_\Delta}{2} \right) \left[+ \log \left(\frac{n_o}{n_\Delta} \right) \right]$$

- the term in square brackets is zero for $n_o = n_\Delta$

Linear Discriminant Analysis

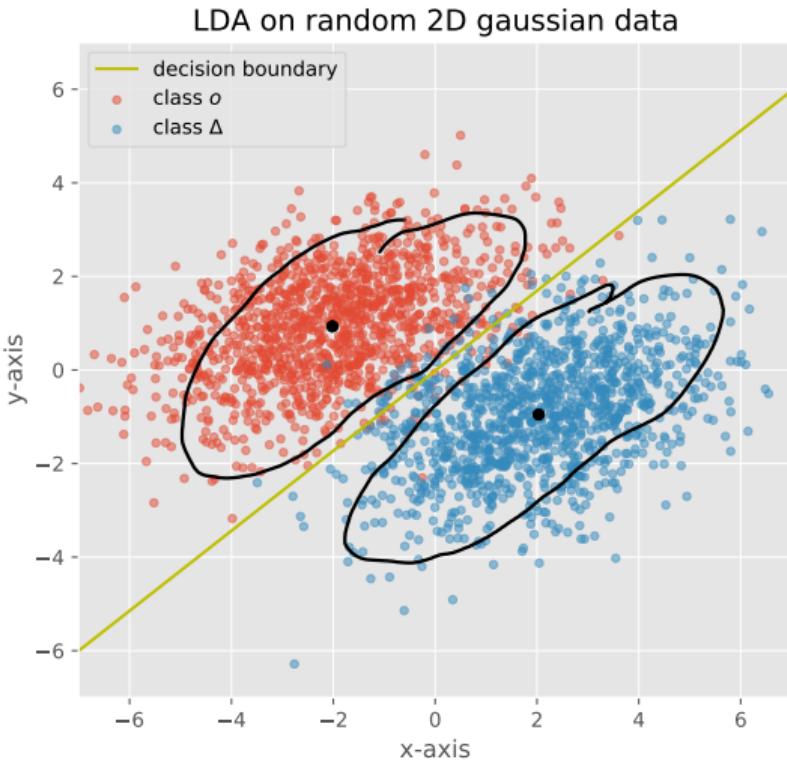
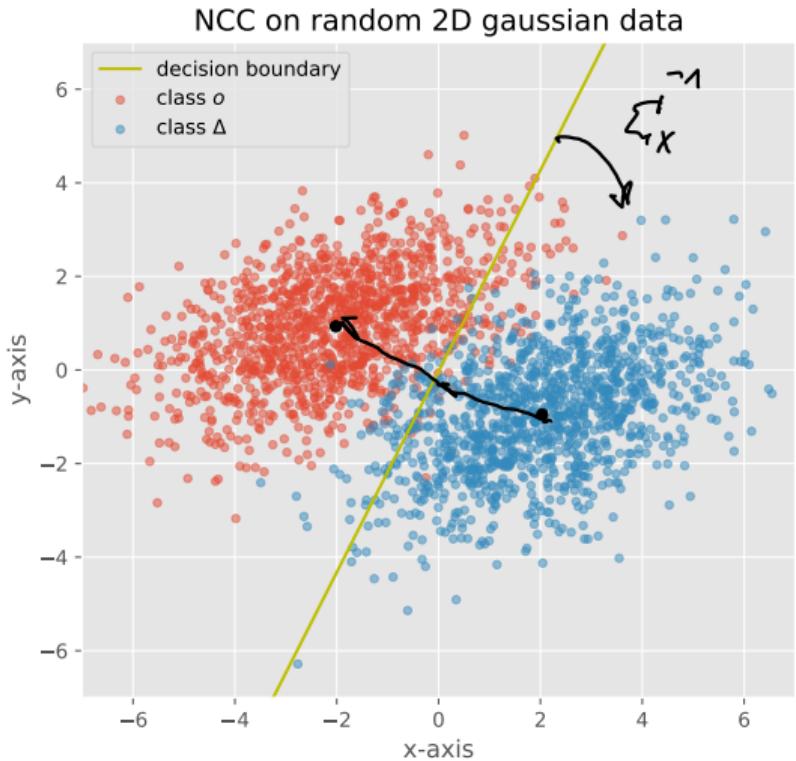
- given class means $\mathbf{w}_o, \mathbf{w}_\Delta \in \mathbb{R}^d$ and number of points n_o, n_Δ per class

$$\mathbf{w} = \Sigma_X^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta)$$

$$\beta = \mathbf{w}^T \left(\frac{\mathbf{w}_o + \mathbf{w}_\Delta}{2} \right) \left[+ \log \left(\frac{n_o}{n_\Delta} \right) \right]$$

- the term in square brackets is zero for $n_o = n_\Delta$
- LDA is the optimal classifier when
 - both classes are gaussian distributed
 - both covariance matrices Σ_X are equal and known

Linear Discriminant Analysis



Task 3 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- ① How does the covariance matrix look like? Which class of matrix is that?



$$\Sigma_K = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \text{diagonal}$$

Task 3 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- ① How does the covariance matrix look like? Which class of matrix is that?

$$\Sigma_X = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \text{this is a diagonal matrix}$$

Task 3 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- ① How does the covariance matrix look like? Which class of matrix is that?

$$\Sigma_X = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \text{this is a diagonal matrix}$$

- ② Which property is represented by the diagonal elements of the covariance matrix?

Task 3 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- ① How does the covariance matrix look like? Which class of matrix is that?

$$\Sigma_X = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \text{this is a diagonal matrix}$$

- ② Which property is represented by the diagonal elements of the covariance matrix?

- σ_1^2 is the variance along the x-axis
- σ_2^2 is the variance along the y-axis

$$\begin{bmatrix} 5 & 0 \\ 0 & 1 \end{bmatrix}$$



Task 3 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- ① How does the covariance matrix look like? Which class of matrix is that?

$$\Sigma_X = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix} \quad \text{this is a diagonal matrix}$$

- ② Which property is represented by the diagonal elements of the covariance matrix?
 - σ_1^2 is the variance along the x-axis
 - σ_2^2 is the variance along the y-axis
 - Note that not necessarily $\sigma_1^2 = \sigma_2^2$

Task 3.3 - NCC vs. LDA

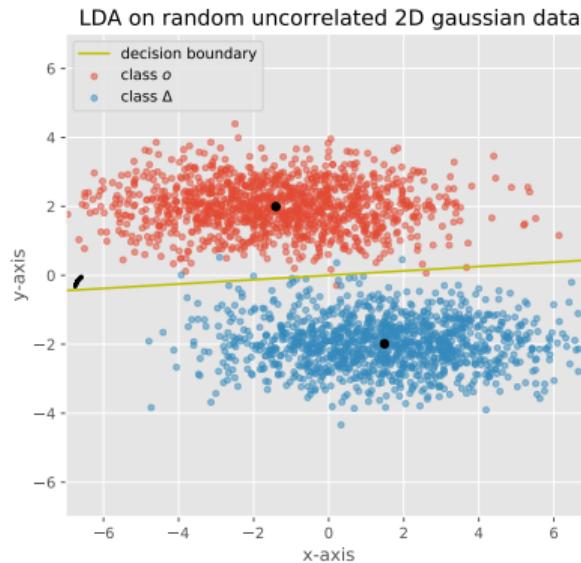
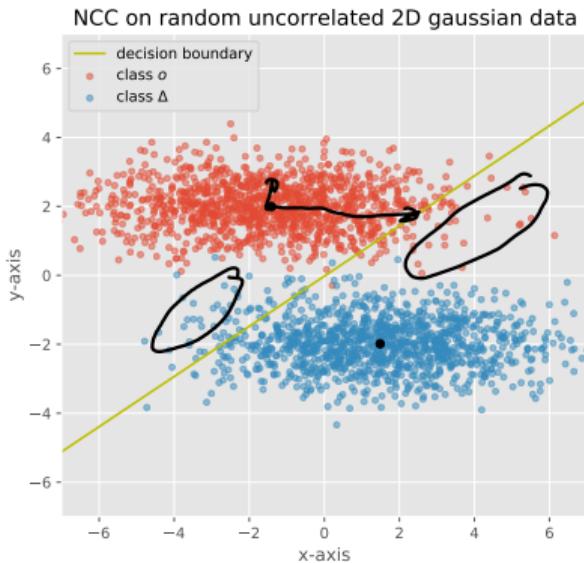
Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- Make a sketch of an imaginary data set with uncorrelated features in which NCC and LDA would classify some data points differently.

Task 3.3 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- Make a sketch of an imaginary data set with uncorrelated features in which NCC and LDA would classify some data points differently.



Task 3.4 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- In which case are NCC and LDA equivalent? Proof that they are equivalent in this case.

Task 3.4 - NCC vs. LDA

Imagine a classification task with two equally gaussian distributed but uncorrelated classes

- In which case are NCC and LDA equivalent? Proof that they are equivalent in this case.

$$\Sigma_X = \alpha \begin{bmatrix} 1 & 0 \\ \cdot & \cdot \\ 0 & 1 \end{bmatrix}, \quad \alpha \in \mathbb{R}$$

Task 3.4 - NCC vs. LDA

$$0 < \omega^T x - \frac{1}{2} \omega^T (\omega_0 + \omega_A) \Leftrightarrow$$

$$0 < (\Sigma_X^{-1}(\omega_0 - \omega_A))^T x - \frac{1}{2} (\Sigma_X^{-1}(\omega_0 - \omega_A))^T (\omega_0 + \omega_A) \Leftrightarrow$$

$$0 < \left(\frac{1}{2} I (\omega_0 - \omega_A) \right)^T x - \frac{1}{2} \left(\frac{1}{2} I (\omega_0 - \omega_A) \right)^T (\omega_0 + \omega_A) \Leftrightarrow$$

$$0 < \frac{1}{2} \left(I (\omega_0 - \omega_A) \right)^T x - \frac{1}{2} \left(I (\omega_0 - \omega_A) \right)^T (\omega_0 + \omega_A) \Leftrightarrow$$

$$0 < \frac{1}{2} \underbrace{(\omega_0 - \omega_A)^T x}_{\omega_{NCC}} - \frac{1}{2} \underbrace{(\omega_0 - \omega_A)^T (\omega_0 + \omega_A)}_{\omega_{NCC}} \quad | \cdot \lambda \quad \begin{cases} \lambda \geq 0 \\ \lambda < 0 \end{cases}$$

$$0 < \omega_{NCC}^T x - \beta_{NCC}$$

$$\left(\Sigma_X^{-1} = \lambda I \right) \Leftrightarrow \Sigma_X^{-1} = \frac{1}{\lambda} I$$

$$\begin{aligned} \Sigma_X^{-1} \Sigma_X^{-1} &= (\lambda I) \left(\frac{1}{\lambda} I \right) \\ &= \lambda \frac{1}{\lambda} (I I) \\ &= I \end{aligned}$$

$$-1 \cdot \omega = \tilde{\omega}$$



- 1 Task 1 - Linear Classification with NCC and Perceptron
- 2 Task 2 - Covariance Matrices
- 3 Task 3 - Linear Discriminant Analysis
- 4 Task 4 - Understanding LDA: the whitening operation
- 5 Model evaluation

Orthogonal Matrices

- matrices with orthonormal column vectors are called orthogonal matrices

$$U = [u_1, \dots, u_d] \quad u_i^T u_j = \begin{cases} 0, & i \neq j \\ 1, & i = j \end{cases}$$
$$\|u_i\| = \|u_i\|^2 = 1$$

Orthogonal Matrices

- matrices with orthonormal column vectors are called orthogonal matrices
- they always have a full rank, thus they are invertible

Orthogonal Matrices

- matrices with orthonormal column vectors are called orthogonal matrices
- they always have a full rank, thus they are invertible
- they are numerical stable when applying them to vectors or other matrices

Orthogonal Matrices

- matrices with orthonormal column vectors are called orthogonal matrices
- they always have a full rank, thus they are invertible
- they are numerical stable when applying them to vectors or other matrices
- for a orthogonal matrix the following equations holds

$$U^T = U^{-1}$$

$$U^T U = U U^T = I$$

Orthogonal Matrices

- matrices with orthonormal column vectors are called orthogonal matrices
- they always have a full rank, thus they are invertible
- they are numerical stable when applying them to vectors or other matrices
- for a orthogonal matrix the following equations holds

$$U^T = U^{-1}$$

$$U^T U = U U^T = I$$

- they not always are symmetric

$$\begin{matrix} V^T = V^{-1} \\ V = V^T \end{matrix} \quad VV = I$$

Orthogonal Matrices

- matrices with orthonormal column vectors are called orthogonal matrices
- they always have a full rank, thus they are invertible
- they are numerical stable when applying them to vectors or other matrices
- for a orthogonal matrix the following equations holds

$$U^T = U^{-1}$$

$$U^T U = UU^T = I$$

- they not always are symmetric
- rotation matrices R are orthogonal matrices with $\det|R| = +1$

Eigendecomposition

- symmetric matrices $A \in \mathbb{R}^{n \times n}$ have n orthogonal eigenvectors and n real eigenvalues

Eigendecomposition

- symmetric matrices $A \in \mathbb{R}^{n \times n}$ have n orthogonal eigenvectors and n real eigenvalues
- thus they can be decomposed into

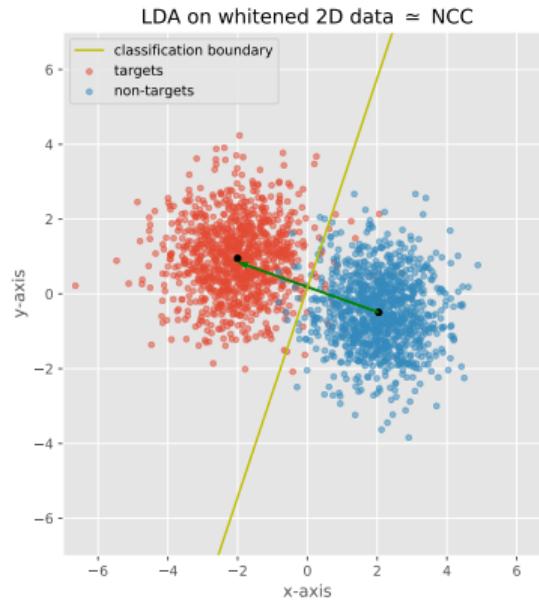
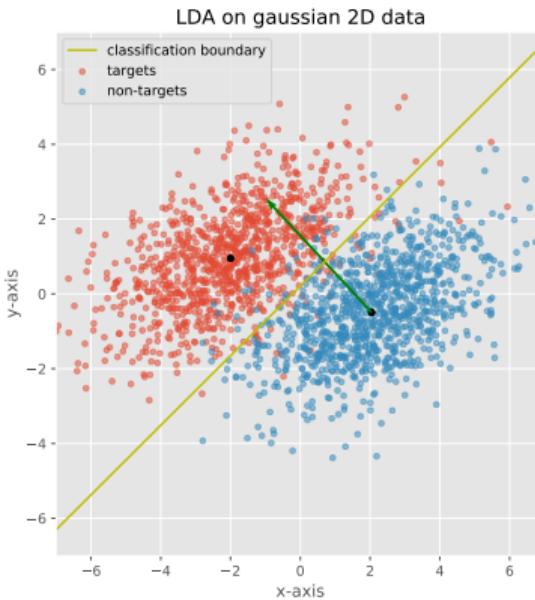
$$A = U\Lambda U^T$$

- U being a orthogonal matrix with the eigenvectors in the columns
- Λ being a diagonal matrix with the eigenvalues at the diagonal

LDA as whitening followed by NCC

LDA first whitens the data and then applies the NCC

$$\sum_{\mathbf{x}}^{\mathbf{I}} = \mathbf{I}$$



LDA as whitening followed by NCC

LDA first whitens the data and then applies the NCC

$$\left(\Sigma_X^{-1} \right)^T = \left(\Sigma_X^T \right)^{-1} \\ = \Sigma^{-1}$$

$$\begin{aligned} \mathbf{w}^T \mathbf{x} &= (\Sigma_X^{-1}(\mathbf{w}_o - \mathbf{w}_\Delta))^T \mathbf{x} \\ &= (\mathbf{w}_o - \mathbf{w}_\Delta)^T \Sigma_X^{-1} \mathbf{x} \quad \downarrow \\ &= (\mathbf{w}_o - \mathbf{w}_\Delta)^T U \Lambda^{-1} U^T \mathbf{x} \\ &= (\mathbf{w}_o - \mathbf{w}_\Delta)^T U \Lambda^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} U^T \mathbf{x} \\ &= \underbrace{(\mathbf{w}_o - \mathbf{w}_\Delta)^T U \Lambda^{-\frac{1}{2}}}_{\text{mean class difference}} \underbrace{\Lambda^{-\frac{1}{2}} U^T \mathbf{x}}_{\substack{\text{whitened } \mathbf{x} \\ \text{of whitened data}}} \end{aligned}$$

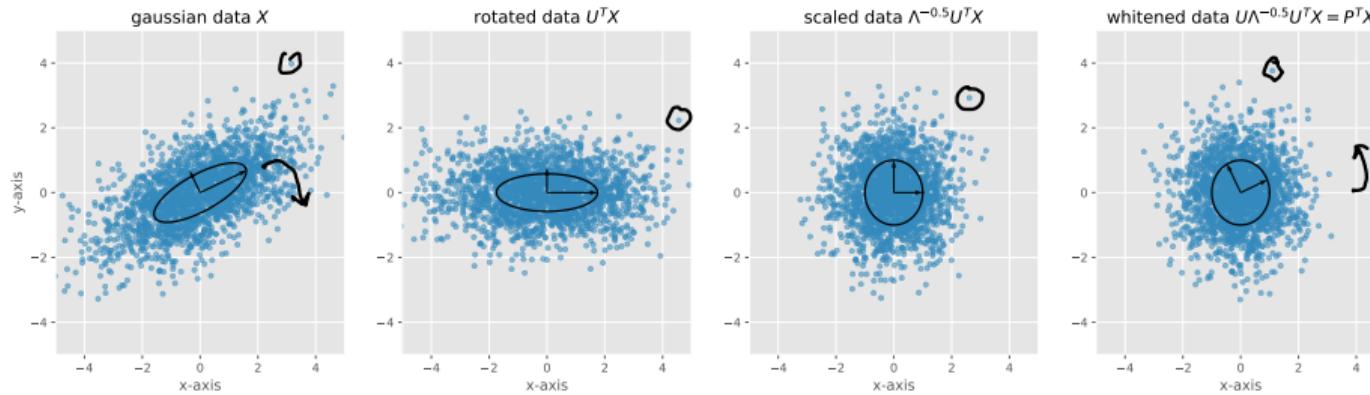
$$\mathbf{w}_{NCC}^T \tilde{\mathbf{x}}$$

Whitening as a "standalone" transformation

In the following tasks we define the mapping $P : \mathbb{R}^d \rightarrow \mathbb{R}^d$ as follows

$$\begin{aligned} P &= U \Lambda^{-\frac{1}{2}} U^\top \\ &= \Sigma_X^{-\frac{1}{2}} \end{aligned}$$

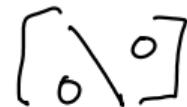
$$\underline{z_k} := P^\top \mathbf{x}_k = \underline{U \Lambda^{-\frac{1}{2}} U^\top \mathbf{x}_k} = \Sigma_X^{-\frac{1}{2}} \mathbf{x}_k$$



Task 4.1 - P is symmetric

Show that $P = U\Lambda^{-\frac{1}{2}}U^T$ is a symmetric matrix

$$\begin{aligned} P &= P^T \quad P^T = (U \Delta^{-\frac{1}{2}} U^T)^T \\ &\quad \text{brace under } \Delta^{-\frac{1}{2}} \\ &= U^T \Delta^{-\frac{1}{2}} U^T \\ &= U \Delta^{-\frac{1}{2}} U^T \\ &= U \Delta^{\frac{1}{2}} U^T = P \end{aligned}$$



Task 4.2 - Matrix square roots of positive definite matrices

Show that $\Sigma_{X^2}^{1/2} = U \Lambda^{1/2} U^T$ is a valid square root of a positive definite matrix

$$\Sigma_{X^2}^{1/2} = U \Lambda^{1/2} U^T \leftarrow$$

$$\begin{aligned}\Sigma_{X^2}^{1/2} \Sigma_{X^2}^{1/2} &= \Sigma_X \\ \Sigma_X^{1/2} \Sigma_X^{1/2} &= (U \Lambda^{1/2} U^T) \underbrace{(U \Lambda^{1/2} U^T)}_{= U \Lambda^{1/2} U^T} \\ &= U \Lambda^{1/2} U^T\end{aligned}$$

$$= U \begin{bmatrix} \lambda_1^{1/2} & & \\ & \ddots & 0 \\ 0 & & \lambda_d^{1/2} \end{bmatrix} \begin{bmatrix} \lambda_1^{1/2} & & \\ & \ddots & 0 \\ 0 & & \lambda_d^{1/2} \end{bmatrix} U^T$$

$$= U \Lambda^{1/2} U^T = \Sigma_X^{1/2}$$

$$U^T U = I$$

Task 4.3 - Inverse of matrix square root

Show that $P = U\Lambda^{-\frac{1}{2}}U^T$ is a valid inverse of the square root of a positive definite matrix

$$\begin{aligned}
 \Sigma_X^{-1} &= U\Lambda^{-1}U^T = PP = U\Lambda^{-\frac{1}{2}}U^T U\Lambda^{-\frac{1}{2}}U^T \\
 &= U\Lambda^{-\frac{1}{2}}\Lambda^{-\frac{1}{2}}U^T \\
 &= U\Lambda^{-1}U^T \\
 &= \Sigma_X^{-1}
 \end{aligned}$$

Task 4.4 - Covariance Matrix of whitened data is the identity

Show that the covariance matrix of whitened data $P^T X$ is the identity Matrix.

Assume that the data has zero mean.

$$\begin{aligned}
 Z &= P^T X \quad \sum_Z = \sum_{I_P^T X} = \frac{1}{n} (P^T X) (P^T X)^T \\
 &= \frac{1}{n} P^T \underline{X X^T} P = P^T \left(\frac{1}{n} X X^T \right) P \\
 &= P^T (\Sigma_X) P \\
 &= P^T U \Lambda U^T P = P U \Lambda U^T P \\
 &= U \Lambda^{\frac{1}{2}} U^T \Lambda \Lambda^{\frac{1}{2}} U^T U \Lambda^{-\frac{1}{2}} U^T \\
 &= U \Lambda^{\frac{1}{2}} \Delta \Delta^{\frac{1}{2}} U^T = U \Lambda^{-\frac{1}{2}} \Lambda^{-\frac{1}{2}} \Delta U^T \\
 &= U \Lambda^{-1} \Delta U^T = U U^T = I
 \end{aligned}$$

$$\Sigma_X = \frac{1}{n} X X^T$$

Task 4.5 - Equivalence of NCC and LDA on whitened data

Proof that classification with LDA is equivalent to classification with NCC on whitened data

$$\begin{aligned}
 \omega_{NCC}^T z &= (\rho^T(\omega_0 - \rho^T \omega_A))^T \rho^T x = (\rho^T(\omega_0 - \omega_A))^T \rho^T x \\
 &= (\omega_0 - \omega_A)^T \rho \rho^T x \\
 &\stackrel{\rho = \Sigma_x^{-\frac{1}{2}}}{=} (\omega_0 - \omega_A)^T \Sigma_x^{-1} x \\
 &\stackrel{\rho \rho = \Sigma_x^{-\frac{1}{2}} \Sigma_x^{-\frac{1}{2}} = \Sigma_x^{-1}}{=} (\underbrace{\Sigma_x^{-1}(\omega_0 - \omega_A)}_{\omega_{LDA}})^T x
 \end{aligned}$$

Task 4.5 - Equivalence of NCC and LDA on whitened data

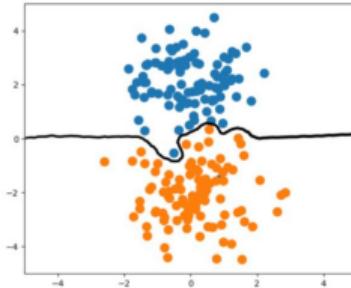
Proof that classification with LDA is equivalent to classification with NCC on whitened data

Generalization

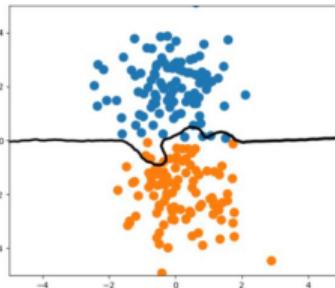
- goal of classification is generalization
- correct categorization of new, unseen data
- correct classification on training data not implies that a classifier is optimal

Overfitting

Performance on
training data

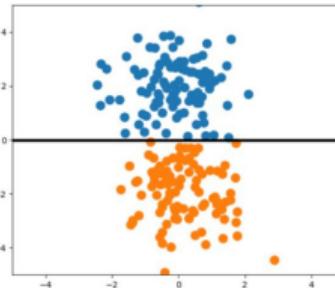
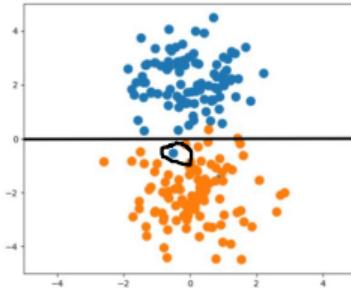


Performance on
new data



overfit

- always compare test accuracy on new data with training accuracy
- should be as similar as possible



Q&A

- The videos are pretty long. I am sorry for that!
- I hope I didn't was to fast or slow. I am happy about feedback :)
- next office hour for questions: wednesday, may 13 at 4pm on zoom