

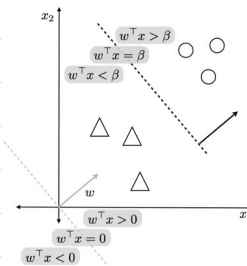
NCC and Perceptron

	NCC	Perceptron
Problem	Classification	Classification, (Regression)
Model	$y = \text{sign}(\mathbf{w}^T \mathbf{x})$	$y = f(\mathbf{w}^T \mathbf{x})$
Error	distance to $\mathbf{w}_{+1}, \mathbf{w}_{-1}$	$-\sum_{m \in M} \mathbf{w}^T \mathbf{x}_m y_m$
Optimization	closed form	SGD
Result	always the same	can differ
Application	Cancer Prediction ¹	NLP ²

NCC:

$$0 < \underbrace{(\mathbf{w}_o - \mathbf{w}_\Delta)^T \mathbf{x}}_{\mathbf{w}} - \underbrace{\frac{1}{2}(\mathbf{w}_o^T \mathbf{w}_o - \mathbf{w}_\Delta^T \mathbf{w}_\Delta)}_{\beta}$$

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



$$\mathbf{x}, \mathbf{w} \in \mathbb{R}^D, \beta \in \mathbb{R}$$

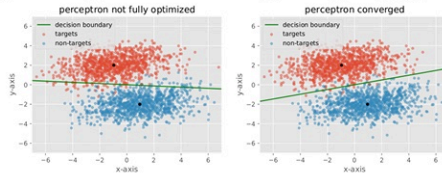
$$\mathbf{w}^T \mathbf{x} - \beta = \begin{cases} > 0 & \text{if } \mathbf{x} \text{ belongs to } o \\ < 0 & \text{if } \mathbf{x} \text{ belongs to } \Delta \end{cases}$$

Points on the decision boundary satisfy $y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} - \beta = 0$

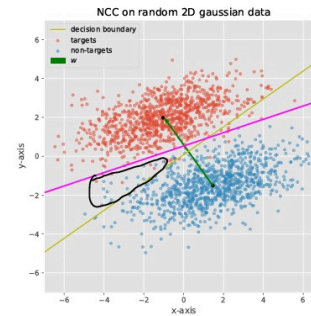
Repetition Perceptron

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

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



Repetition Nearest Centroid Classifier



- compare distance of data points \mathbf{x} to class means
- decide for class with closer center

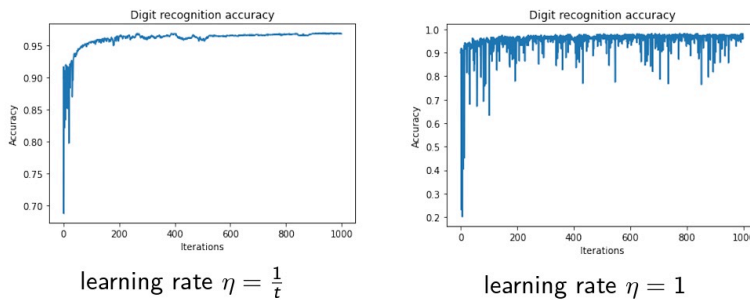
$$\mathbf{w} = \mathbf{w}_o - \mathbf{w}_\Delta$$

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

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

- correlation not considered, thus often not optimal

Perceptron: A good learning rate is useful, but not always necessary to reach convergence.



Task 4 - NCC vs. LDA [6 points]

Last lecture we looked at NCC. This lecture we looked at LDA. These models are closely related, as you saw in the lecture. They are both linear classifiers, which means they both work using a weight vector \mathbf{w} and a bias β , and classify points \mathbf{x}_i using the expression: $\mathbf{w}^T \mathbf{x}_i - \beta \geq 0$. Recall the definition for LDA given in the lecture, given a dataset $X \in \mathbb{R}^{d \times n}$ and class means $\bar{\mathbf{x}}_+, \bar{\mathbf{x}}_- \in \mathbb{R}^d$

$$\mathbf{w} = \bar{S}^{-1}(\bar{\mathbf{x}}_+ - \bar{\mathbf{x}}_-)$$

$$\beta = \frac{1}{2} \mathbf{w}^T (\bar{\mathbf{x}}_+ + \bar{\mathbf{x}}_-) \left[+ \log \left(\frac{n_-}{n_+} \right) \right]$$

vs for NCC

$$\mathbf{w} = \bar{\mathbf{x}}_+ - \bar{\mathbf{x}}_-$$

$$\beta = \frac{1}{2} \mathbf{w}^T (\bar{\mathbf{x}}_+ + \bar{\mathbf{x}}_-)$$

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

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

$$\Sigma_X = \frac{1}{n} (\mathbf{X} - \bar{\mathbf{X}})(\mathbf{X} - \bar{\mathbf{X}})^T$$

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

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

- 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

Goal

Find $\mathbf{w} \in \mathbb{R}^d$ that

- maximizes mean class difference
- minimizes variance in each class

Formalization

Maximize the **Fisher criterion**

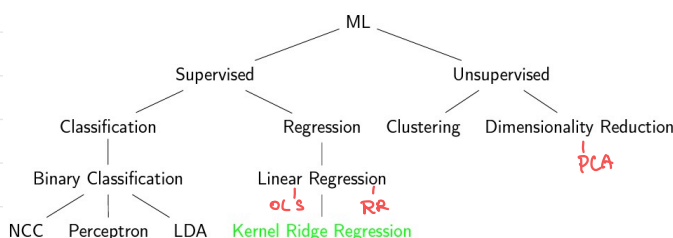
$$J(\mathbf{w}) = \frac{\text{between class variance}}{\text{within class variance}} = \frac{(\mu_o - \mu_\Delta)^2}{\sigma_o^2 + \sigma_\Delta^2}$$

Solution

After some calculations...

$$\mathbf{w} \propto \Sigma_W^{-1}(\bar{\mathbf{x}}_o - \bar{\mathbf{x}}_\Delta)$$

The Tree of CA



S_B : between class scatter

S_W : within class scatter

Linear Regression - The Solution

$$d=1 \rightarrow w = \frac{\sum k_i y_i}{\sum k_i x_i}$$

OLS

- data $x_1, \dots, x_n \in \mathbb{R}^1$ with respective class labels $y_1, \dots, y_n \in \mathbb{R}$
- the function $f(x) = w^T x$ is optimized for

$$\rightarrow w = (XX^T)^{-1} X y^T$$

- in the sense of ordinary least squares, that is

$$\mathcal{E}_{lsq}(w) = \sum_{i=1}^N (y_i - w^T x_i)^2 = \|y - w^T X\|^2$$

- Optimizing $\mathcal{E}_{lsq}(w)$ by setting $\frac{\partial \mathcal{E}_{lsq}(w)}{\partial w} = 0$ leads to the same result.

Ridge Regression

- Regression with regularization: restrict large values for w
- Often it is important to control the complexity of the solution w

$$\mathcal{E}_{RR}(w) = \underbrace{\|y - w^T X\|^2}_{OLS} + \underbrace{\lambda \|w\|^2}_{RR} \quad \frac{d\mathcal{E}_{RR}}{dw} = 0$$

- Solution is given by

biased estimator
but smaller variance

$$w = (XX^T + \lambda I)^{-1} X y^T$$

Hypameter

Underfitting occurs when a model is too simple to capture the underlying structure of the data, while overfitting happens when a model is too complex, capturing noise in the data as if it were signal.

Polynomial Regression - LSE

- This leads again to a LSE, that can be solved

$$X = [\phi(x_1), \phi(x_2), \dots, \phi(x_n)] = \begin{bmatrix} 1 & 1 \\ x_1 & x_2 \\ x_1^2 & x_2^2 \end{bmatrix}$$

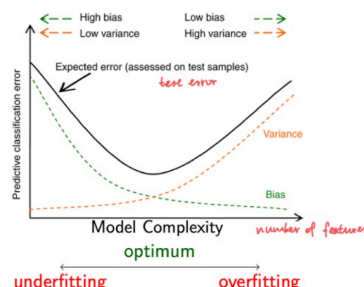
$$\Phi^T w = y^T$$

$$\Phi \Phi^T w = \Phi y^T$$

$$w = (\Phi \Phi^T)^{-1} \Phi y^T$$

- Other choices of basis expansion ϕ yield e.g. spline regression

The bias-variance trade-off



Careful...

Bias and variance are terms with multiple (related) usages

- $w^T x + b$, $w \cdot x$
- Bias/variance of estimator
- Bias/variance of general ML model

Constructing Valid Kernels

1. Option

- Find it by starting with a feature space mapping:
 $k(x, x') = \phi(x)^T \phi(x')$

2. Option

- Show the kernel function corresponds to a scalar product in some feature space
- Option 1: show $K := [k(x_i, x_j)]_{ij}$ is symmetric PSD for all possible choices of $X \subseteq \mathcal{X}$
- Option 2: construct directly from other valid kernels!

Kernelizing Ridge Regression

- Instead we compute:

$$f(x) = w^T \phi(x) = \phi(x)^T w = \underbrace{\phi(x)^T \phi(X)}_{k(x, X)} (K + \lambda I)^{-1} y^T$$

$$= k(x, X) \underbrace{(K + \lambda I)^{-1} y^T}_{\alpha}$$

Instead of w we compute α . What is α ?

$$f(x) = \sum_{i=1}^n \alpha_i k(x, x_i)$$

Kernelized RR - The Algorithm

1. Compute α :

$$\alpha = (K + \lambda I)^{-1} y_{train}^T$$

with $K := k(X_{train}, X_{train})$

2. Predictions:

$$\hat{y}_{new} = \alpha^T k(X_{train}, x_{new})$$

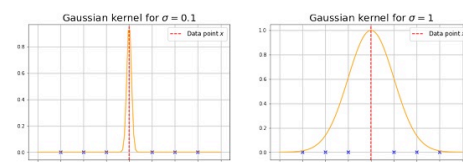
Kernels:

- Unknown (data in) feature space
- $D \gg n$
- Non-numeric data
- The entire dataset has to be stored

Feature maps:

- Feature map often infeasible to compute
- $n \gg D$

Gaussian Kernel



$$k(x, x') := \exp\left\{-\frac{\|x - x'\|^2}{2\sigma^2}\right\}$$

kRF

Repetition: Eigenvalues and Eigenvectors

Given a Matrix $A \in \mathbb{R}^{d \times d}$, then a non-zero vector $\mathbf{v} \in \mathbb{C}^d \setminus \{\mathbf{0}\}$ is called an eigenvector, if there is an eigenvalue $\lambda \in \mathbb{C}$, such that

$$A\mathbf{v} = \lambda\mathbf{v}$$

- eigenvectors are special directions, for which a matrix A is only scales, but not rotate

- A matrix is singular if one or more eigenvalues are zero

$$\lambda = 0 \Leftrightarrow \text{Kern}(A) \neq \emptyset$$

Rotation Matrices II

- Rotation matrices R_θ are always orthogonal with $\det |R_\theta| = +1$

- Are all orthogonal matrices also rotation matrices? (Task 4.2)

- For a orthogonal matrix V the determinant is $\det |V| = \pm 1$

$$V = \begin{pmatrix} 0 & 1 \\ 1 & 0 \end{pmatrix}$$

- V is orthogonal, since $V^T V = I$, but $\det |V| = -1$ and thus it is not a rotation matrix

Summary: Principal Component Analysis

- 1 Estimate the covariance matrix S of the data $X \in \mathbb{R}^{d \times n}$
- 2 Compute the eigenvectors of S
- 3 $W = [\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_k] \in \mathbb{R}^{d \times k}$ where $\mathbf{w}_1, \dots, \mathbf{w}_k \in \mathbb{R}^d$ are the eigenvectors corresponding to the k largest eigenvalues
- 4 Project the data onto W : $H = W^T \cdot X$
- 5 If needed: reconstruct data by $X \approx \tilde{X} = WH$.
This holds for all matrix factorization methods.

Kernel PCA Algorithm

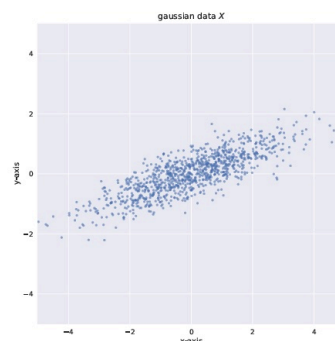
Algorithm 2: Kernel PCA

Input: Dataset $X \in \mathbb{R}^{d \times n}$ with $d \gg n$; Number of PCs k

Output: First k PCs $W \in \mathbb{R}^{d \times k}$, Hidden causes $H \in \mathbb{R}^{k \times n}$ of the dataset

- 1 Compute Kernel $K = (X - \bar{X})^T (X - \bar{X}) \in \mathbb{R}^{n \times n}$
- 2 Compute EVD $K = V\Lambda V^T$
- 3 Take first k eigenvectors corresponding to largest eigenvalues $\alpha = [\mathbf{v}_1, \dots, \mathbf{v}_k]$
- 4 Calculate $W = X\alpha \in \mathbb{R}^{d \times k}$
- 5 Project data $H = W^T X$
- 6 **return** W and H

Principle Component Analysis

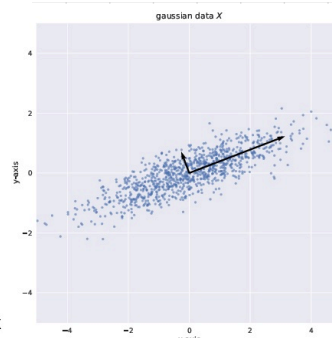


- Which dimensions contain the most information?

$$\Rightarrow \text{directions } \mathbf{w} \text{ with the highest variance} \\ \arg\max_{\mathbf{w}} \text{var}(\mathbf{w}^T X) = \arg\max_{\mathbf{w}} \mathbf{w}^T \Sigma_X \mathbf{w}$$

- directions \mathbf{w} that minimizes the noise

$$\arg\min_{\mathbf{w}} \sum_{i=1}^n \|(\mathbf{w}^T \mathbf{x}_i)\mathbf{w} - \mathbf{x}_i\|^2 = \arg\max_{\mathbf{w}} \mathbf{w}^T \Sigma_X \mathbf{w}$$



- $\mathbf{w}^T \Sigma_X \mathbf{w}$ is maximized by the eigenvector of Σ_X with the largest eigenvalue
- The eigenvectors of Σ_X are the principle components

- PCs can be constrained to unit length $\|\mathbf{w}\| = 1$

- Variance along PC is given by the corresponding eigenvalue, since

$$\text{var}(\mathbf{w}^T X) = \mathbf{w}^T \Sigma_X \mathbf{w} = \mathbf{w}^T \lambda \mathbf{w} = \lambda \mathbf{w}^T \mathbf{w} = \lambda$$

Kernel PCA

- If dimensionality becomes larger than number of samples $d \gg n$
 - Only $\leq n$ non-zero eigenvalues
 - Covariance Σ_X becomes singular, $\text{Rank}(\Sigma_X) \leq n$
 - Covariance is positive *semi* definite
 - Covariance $\Sigma_X \in \mathbb{R}^{d \times d}$ will be very large
 - High time and space complexity for calculation

- Kernel PCA should be used to estimate n first Principal Components

From PCA to NMF

- Given data $X \in \mathbb{R}^{d \times n}$ and PCs in the columns of a matrix $W \in \mathbb{R}^{d \times k}$
- Data can be projected onto PCs $H = W^T X \Leftrightarrow X \approx WH$

- If $X > 0$, it should be possible to choose W, H such that $W > 0, H > 0$

- Goal of NMF is to find $W \in \mathbb{R}^{d \times k}$ and $H \in \mathbb{R}^{k \times n}$, such that $\|X - WH\|_{\text{Fro}}^2$ is minimized
- Cannot be directly calculated \Rightarrow iterative optimization with gradient descent