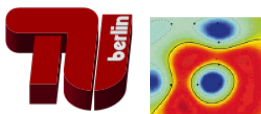


Lecture 0: Introduction

Machine Learning - Theory and Applications



TU Berlin – SoSe 2016

Bayesian Decision Theory

Known:

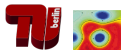
- ▶ Classes $\{\omega_1, \dots, \omega_c\}$
- ▶ Class probabilities $P(\omega_j)$
- ▶ Observation $\mathbf{x} \in \mathbb{R}^d$
- ▶ Likelihood function $P(\mathbf{x}|\omega_j)$

Baye's Rule:

$$p(\omega_j|\mathbf{x}) = \frac{P(\mathbf{x}|\omega_j)P(\omega_j)}{P(\mathbf{x})}$$

Law of total probability:

$$p(\mathbf{x}) = \sum_{j=1}^c p(\mathbf{x}|\omega_j)P(\omega_j)$$

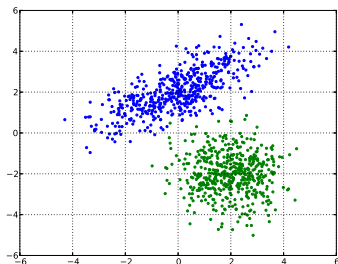


Example of Probability Distribution

Normal distribution

$$P(\mathbf{x}|\omega_j) = \frac{1}{(2\pi)^{\frac{d}{2}} |\boldsymbol{\Sigma}_j|^{\frac{1}{2}}} e^{-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu}_j)^\top \boldsymbol{\Sigma}_j^{-1}(\mathbf{x}-\boldsymbol{\mu}_j)}$$

Parameters $\theta_j = \{\boldsymbol{\mu}_j, \boldsymbol{\Sigma}_j\}$.



Discriminant Functions

Not always interested knowing exactly $p(\omega_j|\mathbf{x})$ for all j .

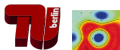
Instead:

$$j^* = \arg \max_j p(\omega_j|\mathbf{x}) = \arg \max_j g_j(\mathbf{x}).$$

where $g_j(\mathbf{x})$ is a discriminant function.

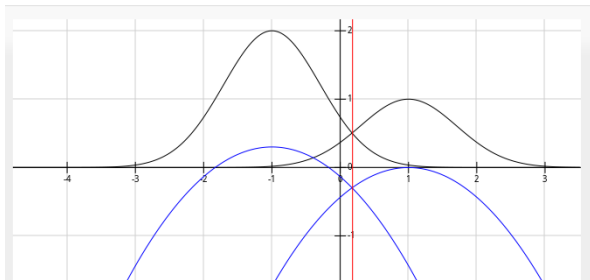
Example of discriminant functions:

- ▶ $g_j(\mathbf{x}) = p(\omega_j|\mathbf{x})$
- ▶ $g_j(\mathbf{x}) = p(\omega_j|\mathbf{x}) + a$
- ▶ $g_j(\mathbf{x}) = p(\omega_j|\mathbf{x}) \cdot \alpha \quad \alpha > 0$
- ▶ $g_j(\mathbf{x}) = \log p(\omega_j|\mathbf{x})$



Discriminant Functions

- ▶ $(g_1(\mathbf{x}), g_2(\mathbf{x}))$
- ▶ $(\log g_1(\mathbf{x}), \log g_2(\mathbf{x}))$



(plotted with www.fooplot.com)



Observation: Both pairs of discriminants meet at the same point on the x-axis (shown by the red line).

Discriminative Learning

Example: Perceptron

- ▶ Labeled dataset:

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N \quad \mathbf{x}_i \in \mathbb{R}^d \quad y_i \in \{1, \dots, c\}$$

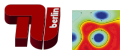
- ▶ Discriminants:

$$g_j(\mathbf{x}) = \mathbf{w}_j^\top \mathbf{x} - b_j \quad j \in \{1, \dots, c\}$$

- ▶ Objective:

$$\max_{\{\mathbf{w}_j, b_j\}} \sum_{i=1}^N \mathbf{1}_{\arg \max_j g_j(\mathbf{x}_i) = y_i}$$

Remaining question: What if several different solutions classify the problem perfectly? Which one to choose?



Model Selection

From which function class should we choose g_j ?

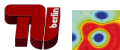
- ▶ **Good idea:** Choose function class with low VC-dimension h .
Example: Large-margin classifiers

$$h \leq \min \left(\left\lceil \frac{R^2}{\Delta^2} \right\rceil, d \right) + 1.$$

where Δ is size of the margin, R is the radius of the minimum enclosing sphere of \mathcal{D} and d is the number of dimensions of the input space.

- ▶ **Bad idea:** Choose function class with small number of parameters. Example: $\sin(ax)$ with $a \in \mathbb{R}$ has only one parameter, but can shatter infinitely many data points on \mathbb{R} .

$$h \leq \infty.$$

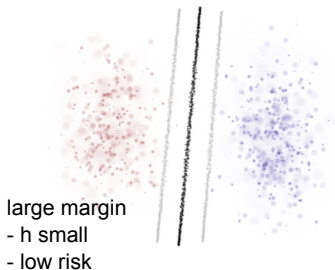
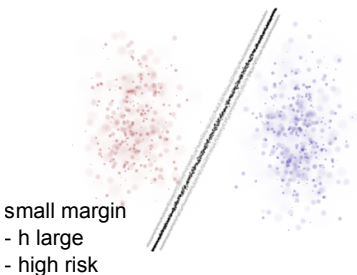


Model Selection

Low VC dimension gives guarantees on generalization error of the model:

$$R[f] \leq R_{\text{emp}}[f] + \sqrt{\frac{h(\log \frac{2N}{h} + 1) - \log(\eta/4)}{N}}$$

with probability $1 - \eta$.



Enforcing Large Margin In Practice

Support Vector Machine

- ▶ Labeled dataset:

$$\mathcal{D} = \{\mathbf{x}_i, y_i\}_{i=1}^N \quad \mathbf{x}_i \in \mathbb{R}^d \quad y_i \in \{-1, 1\}$$

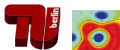
- ▶ Classifier:

$$g_2(\mathbf{x}) - g_1(\mathbf{x}) = \mathbf{w}^\top \mathbf{x} - b$$

- ▶ Objective:

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad (\mathbf{w}^\top \mathbf{x}_i - b) \cdot y_i \geq 1$$

Various techniques to optimize SVMs: quadratic optimization, subgradient methods.



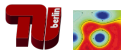
NonLinear Classification via Kernelization

Step 1: Define *Lagrangian*:

$$L(\mathbf{w}, \alpha) = \underbrace{\frac{1}{2} \|\mathbf{w}\|^2}_{\text{objective}} + \sum_{i=1}^n \alpha_i \cdot \underbrace{[1 - (\mathbf{w}^\top \mathbf{x}_i - b) \cdot y_i]}_{(>0) \Rightarrow \text{constraint violation}}$$

Minimize with \mathbf{w}, b under maximally penalized constraints violations ($\max_{\alpha \succeq 0}$):

$$\begin{aligned} \min_{\mathbf{w}, b} \max_{\alpha \succeq 0} L(\mathbf{w}, \alpha) &= \max_{\alpha \succeq 0} \min_{\mathbf{w}, b} L(\mathbf{w}, \alpha) \\ &= \max_{\alpha \succeq 0} \sum_{i=1}^n \alpha_i - \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j \\ &\quad \text{with } \sum_i \alpha_i y_i = 0 \end{aligned}$$



NonLinear Classification via Kernelization

Step 2: Replace $\mathbf{x}_i^\top \mathbf{x}_j$ by a kernel function $k(\mathbf{x}_i, \mathbf{x}_j)$ that must satisfy

$$\sum_{i=1}^N \sum_{j=1}^N c_i c_j k(\mathbf{x}_i, \mathbf{x}_j) \geq 0 \quad \text{for all } c_i, c_j \in \mathbb{R}$$

If $k(\mathbf{x}_i, \mathbf{x}_j)$ is a kernel, then, there exists a feature map $\phi(\mathbf{x})$, possibly infinite-dimensional such that

$$k(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^\top \phi(\mathbf{x}')$$

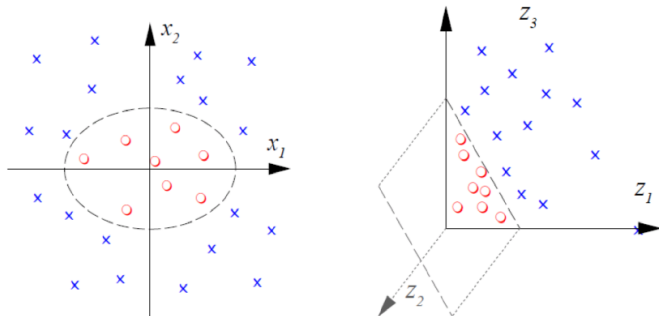
Therefore, a kernel SVM is equivalent to a nonlinear SVM built directly on the nonlinear representation $\phi(\mathbf{x})$.

NonLinear Classification via Kernelization

Example: all second order monomials

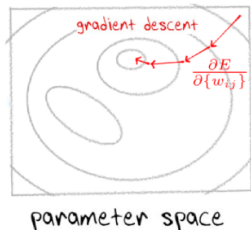
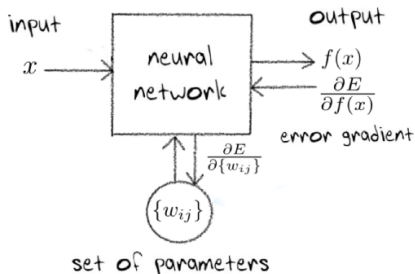
$$\Phi: \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)$$



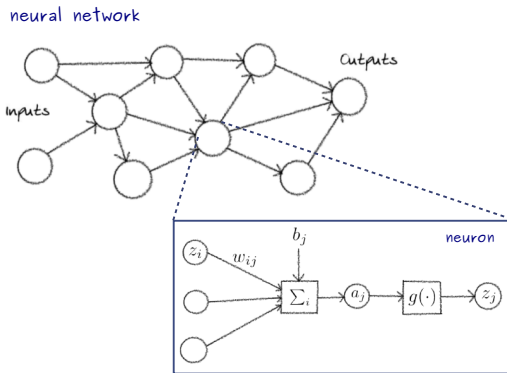
Another Approach to NonLinear Learning

Replace the linear discriminant $g_2(x) - g_1(x) = \mathbf{w}^\top \mathbf{x} - b$ by a more complex nonlinear function composed of many trainable parameters, and for which a gradient can be computed easily. Example: neural networks.



Neural Networks

A neural network is a graph of many simple interconnected computational units (neurons), that together form a complex function.



Gradient is obtained using the error backpropagation algorithm (an efficient application of the chain rule in a graph).