

# Supervised and Unsupervised Learning

## Machine Learning II (SS 2008, TU Berlin)

Prof. Dr. Klaus-Robert Müller  
Dr. Alexander Zien

15. 04. 2008



MAX-PLANCK-GESellschaft

- 1 Overview
  - Taxonomy of Learning Tasks
  - Taxonomy of Learning Approaches
- 2 Regression
- 3 Classification
  - Support Vector Machine
  - Logistic Regression
- 4 Handling Non-Linearity with Kernels
- 5 Spectral Clustering
- 6 Bibliography

# Taxonomy of Learning Tasks

input	output	
	discrete	continuous
<b>supervised</b> given $\{(\mathbf{x}_i, y_i)\}$	<ul style="list-style-type: none"><li>● classification</li></ul>	<ul style="list-style-type: none"><li>● regression</li></ul>
<b>unsupervised</b> given $\{\mathbf{x}_i\}$	<ul style="list-style-type: none"><li>● clustering</li><li>● anomaly detection</li></ul>	<ul style="list-style-type: none"><li>● dimensionality reduction</li></ul>

Outside this scheme:

- active learning
- reinforcement learning

# Taxonomy of Learning Approaches

## Models

- geometrical models
- statistical learning theory
- probabilistic models

## Methods

- optimization
  - regularization
  - ML / MAP
- Bayesian inference

we'll see close connections...

# Regression: Geometric View (1)

Minimize squared distance of predictions to labels:

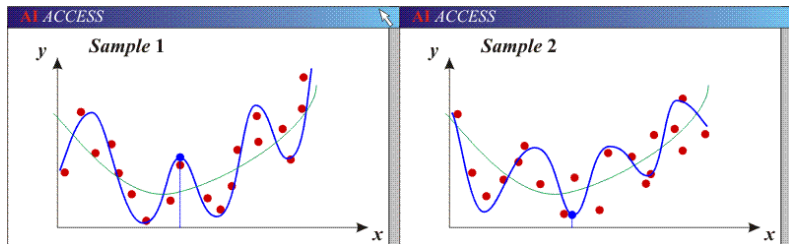
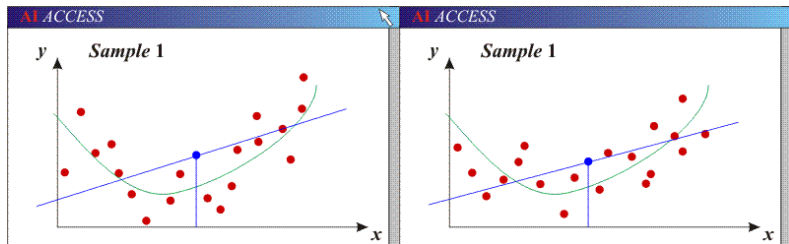
- $J(\mathbf{w}) := \sum_{i=1}^N \left( y_i - \mathbf{x}_i^\top \mathbf{w} \right)^2 = \left( \mathbf{y} - \mathbf{X}^\top \mathbf{w} \right)^\top \left( \mathbf{y} - \mathbf{X}^\top \mathbf{w} \right)$   
where  $\mathbf{X} \in \mathbb{R}^{D \times N}$  contains  $N$  points as columns
- $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} J(\mathbf{w})$

Convex  $\Rightarrow$  vanishing derivative indicates global optimum  $\hat{\mathbf{w}}$ .

$$0 = \left. \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}} \right|_{\mathbf{w}=\hat{\mathbf{w}}} = \mathbf{X}\mathbf{X}^\top \hat{\mathbf{w}} - \mathbf{X}\mathbf{y} \quad \Rightarrow \quad \hat{\mathbf{w}} = \left( \mathbf{X}\mathbf{X}^\top \right)^{-1} \mathbf{X}\mathbf{y}$$

**“Ordinary Least Squares”**

# Regression: Bias-Variance Tradeoff (1)



## Regression: Bias-Variance Tradeoff (2)

Bias-variance decomposition of approximation error:

$$\mathbb{E} \left[ \|\mathbf{y} - \mathbf{X}\hat{\mathbf{w}}\|^2 \right] = \begin{cases} \mathbb{E} \left[ \|\mathbf{y} - \mathbf{X}\mathbf{w}\|^2 \right] & \textit{noise} \\ + \mathbb{E} \left[ \|\mathbf{X}\mathbf{w} - E[\mathbf{X}\hat{\mathbf{w}}]\|^2 \right] & \textit{bias}^2 \\ + \mathbb{E} \left[ \|\mathbb{E}[\mathbf{X}\hat{\mathbf{w}}] - \mathbf{X}\hat{\mathbf{w}}\|^2 \right] & \textit{variance} \end{cases}$$

### Shrinkage

- increase bias...
- ... to reduce variance even more
- first proposed for the multivariate mean (James/Stein)

## Regression: Geometric View (2)

Apply shrinkage to least squares ...

- shrink  $\mathbf{w}$

- $\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\| \mathbf{y} - \mathbf{X}^T \mathbf{w} \right\|^2 + \lambda \|\mathbf{w}\|^2$

- $\Rightarrow \hat{\mathbf{w}} = \left( \mathbf{X}\mathbf{X}^T + \lambda \mathbf{I} \right)^{-1} \mathbf{X}\mathbf{y}$

### “Regularized Least Squares”

- aka ridge regression
- aka Tikhonov regularization



## Regression: Probabilistic View (1)

Multivariate Gaussian error model:

$$P(\mathbf{y}|\mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y} | \mathbf{X}^\top \mathbf{w}, \sigma^2 \mathbf{I})$$

**Likelihood** of parameter  $\mathbf{w}$  given the data  $(\mathbf{X}, \mathbf{y})$ :

$$\mathcal{L}(\mathbf{w}) = (2\pi)^{-\frac{N}{2}} |\sigma^2 \mathbf{I}|^{-\frac{1}{2}} \exp \left[ -\frac{1}{2} (\mathbf{y} - \mathbf{X}^\top \mathbf{w})^\top (\sigma^2 \mathbf{I})^{-1} (\mathbf{y} - \mathbf{X}^\top \mathbf{w}) \right]$$

(same as  $P(\mathbf{y}|\mathbf{X}, \mathbf{w})$ , just seen as function of  $\mathbf{w}$ )

**Maximum Likelihood (ML)** approach:

- $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \log \mathcal{L}(\mathbf{w}) = \arg \min_{\mathbf{w}} \left\| \mathbf{y} - \mathbf{X}^\top \mathbf{w} \right\|^2$
- same as OLS!

## Regression: Probabilistic View (2)

### Maximum A Posteriori (MAP) approach:

- add normal prior:  $P(\mathbf{w}) = \mathcal{N}(\mathbf{0}, \lambda\mathbf{I})$
- maximize posterior  $P(\mathbf{w}|\mathbf{X}, \mathbf{y}) \propto P(\mathbf{w})P(\mathbf{y}|\mathbf{X}, \mathbf{w})$   
(Bayes theorem)
- $\Rightarrow \hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \left\| \mathbf{y} - \mathbf{X}^T \mathbf{w} \right\|^2 + \lambda \|\mathbf{w}\|^2$
- same as RLS!

### Remarks:

- 1 nice way to incorporate prior knowledge:  $P(\mathbf{w}) = \mathcal{N}(\mathbf{w}_0, \lambda\mathbf{I})$
- 2 point estimate — not Bayesian!

## Regression: Probabilistic View (3)

**Bayesian** approach (“inference”):

- follow Bayes' rule:
  - ① encode uncertainty in beliefs as prior distribution
  - ② update beliefs according to data
  - ③ this yields posterior beliefs *as distribution*
- if you believe in 3 simple axioms (Cox/Jaynes), the only way!
- “loss” only considered afterwards ( $\rightarrow$  decision theory)

For our regression setting:

- posterior  $P(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{P(\mathbf{y}|\mathbf{X}, \mathbf{w})P(\mathbf{w})}{P(\mathbf{y}|\mathbf{X})} = \mathcal{N}(\cdot, \cdot)$

- predictive distribution:

$$P(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{y}) = \int P(\mathbf{y}^*|\mathbf{x}^*, \mathbf{w})P(\mathbf{w}|\mathbf{X}, \mathbf{y})d\mathbf{w} = \mathcal{N}(\cdot, \cdot)$$

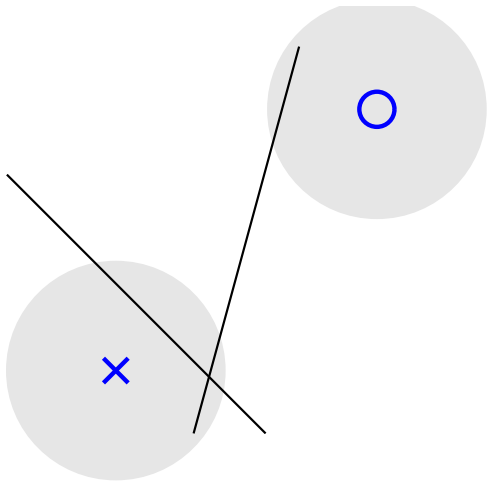
# Classification

## Generative Learning (Sampling Paradigm)

- **model**  $P(y, \mathbf{x})$ , often as  $P(y)P(\mathbf{x}|y)$
- predict via Bayes theorem: 
$$P(y|\mathbf{x}) = \frac{P(y)P(\mathbf{x}|y)}{\sum_{y'} P(y')P(\mathbf{x}|y')}$$
- naive Bayes: assume  $P(\mathbf{x}|y) = \prod_{i=1}^D P(x_{[i]}|y)$
- in general: prior knowledge explicitly built in

## Discriminative Learning (Diagnostic Paradigm)

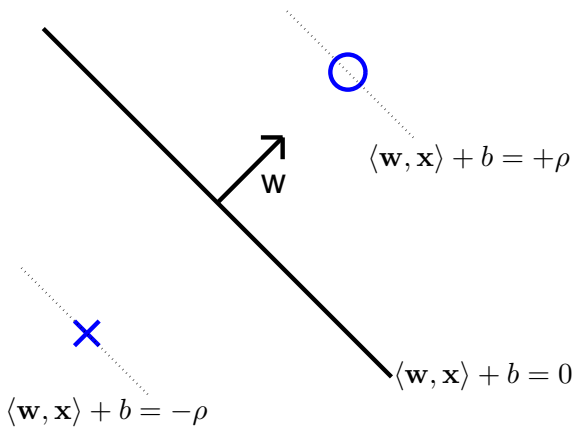
- **model**  $p(y|\mathbf{x})$  (or just boundary:  $\{\mathbf{x} \mid p(y|\mathbf{x}) = \frac{1}{2}\}$ )
- no naive independence assumption
- in general: less prior knowledge (lower bias, higher variance)
- examples: **SVM, Logistic Regression**



not robust wrt input noise!

## SVM:

maximum margin classifier



$$\max_{\mathbf{w}, b, \rho}$$

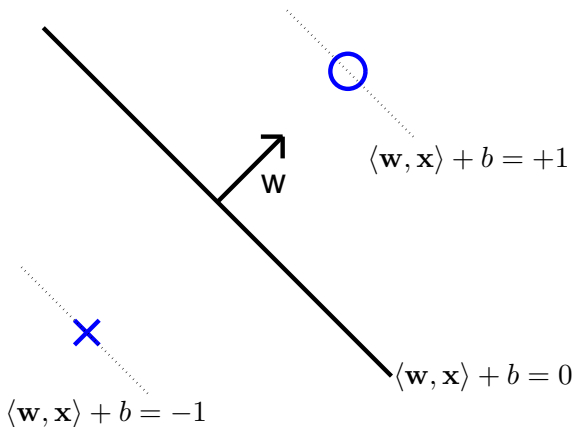
$\underbrace{\rho}_{\text{margin}}$

*s.t.*

$$\underbrace{y_i(\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho}_{\text{data fitting}}$$

$$\underbrace{\|\mathbf{w}\| = 1}_{\text{normalization}}$$

**SVM:**  
regularized  
data fitting



$$\min_{\mathbf{w}, b} \underbrace{\frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle}_{\text{regularizer}} \quad s.t. \quad \underbrace{y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1}_{\text{data fitting}}$$

# Equivalent Reformulation of the SVM

$$\begin{aligned} & \max_{\mathbf{w}, b, \rho} \quad \rho & \text{s.t.} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq \rho, \quad \|\mathbf{w}\| = 1 \\ \Leftrightarrow & \max_{\mathbf{w}', b, \rho} \quad \rho^2 & \text{s.t.} \quad & y_i \left( \left\langle \frac{\mathbf{w}'}{\|\mathbf{w}'\|}, \mathbf{x}_i \right\rangle + b \right) \geq \rho, \quad \rho \geq 0 \\ \Leftrightarrow & \max_{\mathbf{w}', b, \rho} \quad \rho^2 & \text{s.t.} \quad & y_i \left( \underbrace{\left\langle \frac{\mathbf{w}'}{\|\mathbf{w}'\| \rho}, \mathbf{x}_i \right\rangle}_{\mathbf{w}''} + \underbrace{\frac{b}{\rho}}_{b''} \right) \geq 1, \quad \rho \geq 0 \\ \Leftrightarrow & \max_{\mathbf{w}'', b''} \quad \frac{1}{\|\mathbf{w}''\|^2} & \text{s.t.} \quad & y_i (\langle \mathbf{w}'', \mathbf{x}_i \rangle + b'') \geq 1, \end{aligned}$$

$$\text{using } \|\mathbf{w}''\| = \left\| \frac{\mathbf{w}'}{\|\mathbf{w}'\| \rho} \right\| = \left| \frac{1}{\rho} \right| \cdot \left\| \frac{\mathbf{w}'}{\|\mathbf{w}'\|} \right\| = \frac{1}{\rho}$$

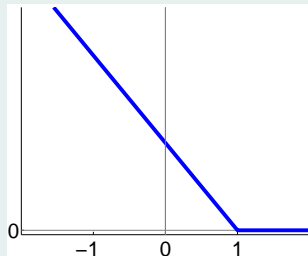


# Soft-Margin SVM Loss

$$\begin{aligned} \min_{\mathbf{w}, b, (\xi_k)} \quad & \frac{1}{2} \langle \mathbf{w}, \mathbf{w} \rangle + C \sum_i \xi_i \\ \text{s.t.} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 - \xi_i, \quad \xi_i \geq 0 \end{aligned}$$

## Effective Loss Function

$$\xi_i = \max \{1 - y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b), 0\}$$



$$y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b)$$

# Logistic Regression (1)

- **log-linear** likelihood ratio:

$$f(\mathbf{x}) := \log \left( \frac{p(y = +1|\mathbf{x})}{p(y = -1|\mathbf{x})} \right) = \mathbf{w}^\top \Phi(\mathbf{x}) + b$$

$$\Rightarrow \text{prediction for } \mathbf{x}: \text{sign} \left( \mathbf{w}^\top \Phi(\mathbf{x}) + b \right)$$

- implied **likelihood**:

$$p(y = +1|\mathbf{x}) = \frac{1}{1 + \exp(-f(\mathbf{x}))}$$

$$p(y = -1|\mathbf{x}) = \frac{1}{1 + \exp(+f(\mathbf{x}))}$$

- possibly **Gaussian prior**

$$\mathbf{w} \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

## Logistic Regression (2)

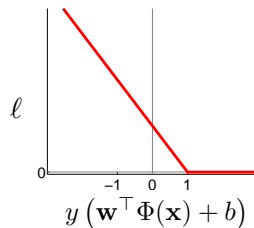
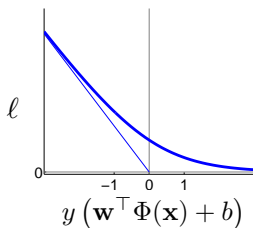
- **maximum likelihood (ML)** estimation: (convex)

$$\min_{\mathbf{w}, b} \sum_i \underbrace{\log \left( 1 + \exp \left( - y_i (\mathbf{w}^\top \Phi(\mathbf{x}_i) + b) \right) \right)}_{=:\ell_{\mathbf{w}, b}(\mathbf{x}_i, y_i)}$$

- **maximum a posteriori (MAP)** estimation:

$$\min_{\mathbf{w}, b} \lambda \|\mathbf{w}\|^2 + \sum_i \ell_{\mathbf{w}, b}(\mathbf{x}_i, y_i)$$

- comparing **LogReg likelihood**  $\ell_{\mathbf{w}, b}$  to **SVM loss**  $\ell_{\mathbf{w}, b}$



# Representer Theorem

Objective:  $J(\mathbf{w}) = \|\mathbf{w}\|^2 + \sum_i \ell_i \left( \mathbf{w}^\top \mathbf{x}_i \right)$  .

## Representer Theorem:

$\mathbf{w}^* := \arg \max_{\mathbf{w}} J(\mathbf{w})$  is in the span of the data  $(\mathbf{x}_i)$ , ie

$$\mathbf{w}^* = \sum_i \alpha_i \mathbf{x}_i .$$

**Proof:** Let  $\mathbf{w}^* = \underbrace{\sum_i \alpha_i \mathbf{x}_i}_{=:\mathbf{w}_{\parallel}} + \mathbf{w}_{\perp}$  with  $\mathbf{w}_{\perp} \perp \mathbf{w}_{\parallel}$ . Then

$$J(\mathbf{w}^*) = \|\mathbf{w}_{\parallel}\|^2 + \|\mathbf{w}_{\perp}\|^2 + \sum_i \ell_i \left( \mathbf{w}_{\parallel}^\top \mathbf{x}_i + \mathbf{w}_{\perp}^\top \mathbf{x}_i \right) = J(\mathbf{w}_{\parallel}) + \|\mathbf{w}_{\perp}\|^2$$



# Non-Linearity via Kernels

## Kernel Functions

Use feature map  $\Phi(\mathbf{x})$ , **kernel**  $k(\mathbf{x}_i, \mathbf{x}_j) = \langle \Phi(\mathbf{x}_i), \Phi(\mathbf{x}_j) \rangle$ .

- Intuitively, kernel measures similarity of two objects  $\mathbf{x}$ .
- Fct is kernel  $\Leftrightarrow$  fct is *positive semi-definite*.

Kernelization possible if data access only through dot products:

- requires  $l_2$ -regularization:  $\|\mathbf{w}\|_2 = \langle \mathbf{w}, \mathbf{w} \rangle$
- SVMs, LogReg, LS-Regression, GPs, ...

# Three Routes to Kernelization

- 1 Dualization (the classic):  
eg SVM:  $\min_{\alpha} \alpha^{\top} \mathbf{H} \alpha - \mathbf{1}^{\top} \alpha$  with  $H_{ij} = y_i y_j k(\mathbf{x}_i, \mathbf{x}_j)$
- 2 Plug in Representer Theorem:  
 $\mathbf{w} = \sum_i \alpha_i \Phi(\mathbf{x}_i)$ ; now optimize  $\alpha$  instead of  $\mathbf{w}$
- 3 Re-represent data:  $E := \text{span}\{\Phi(\mathbf{x}_i)\} \hat{=} \mathbb{R}^N$  (Repr. Thrm.)

1 expand basis vectors  $\mathbf{v}_i$  of  $E$ : 
$$\mathbf{v}_i = \sum_k A_{ik} \Phi(\mathbf{x}_k)$$

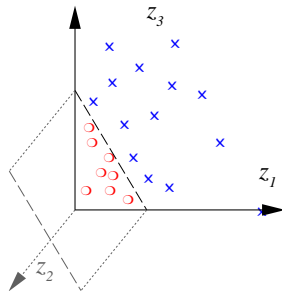
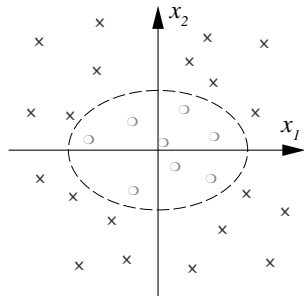
2 orthonormality gives: 
$$(A^{\top} A)^{-1} = K$$
  
solve for  $A$ , eg by KPCA or Choleski decomposition

3 project data  $\Phi(\mathbf{x}_i)$  on basis  $V = (\mathbf{v}_j)_j$ :  
$$\tilde{\mathbf{x}}_i = V^{\top} \Phi(\mathbf{x}_i) = (A)_i$$

# Non-Linear Mappings

## Example: All Degree 2 Monomials

$$\begin{aligned}\Phi : \mathbb{R}^2 &\rightarrow \mathbb{R}^3 =: \mathcal{H} \quad (\text{"Feature Space"}) \\ (x_1, x_2) &\mapsto (z_1, z_2, z_3) := (x_1^2, \sqrt{2} x_1 x_2, x_2^2)\end{aligned}$$



## Example: All Degree 2 Monomials for a 2D Input

$$\begin{aligned}\langle \Phi(x), \Phi(x') \rangle &= \left\langle (x_1^2, \sqrt{2} x_1 x_2, x_2^2), (x_1'^2, \sqrt{2} x_1' x_2', x_2'^2) \right\rangle \\ &= (x_1 x_1' + x_2 x_2')^2 \\ &= \langle x, x' \rangle^2 \\ &=: k(x, x')\end{aligned}$$

→ the dot product in  $\mathcal{H}$  can be computed in  $\mathbb{R}^2$



## Popular Discriminative (Kernel-) Classifiers

method	SVM	Logistic Regression	Fisher Linear Discriminant
models	$p(y x) = 0.5$	$p(y x)$	$p(y x)$
probabilistic	no	yes	yes
coefficients $\alpha$	sparse $\Rightarrow$ efficient optimization	full	full
difference to SVM	—	uses logistic loss fct.	maximizes average margin

# Parametric vs Non-Parametric

Two alternative views (depending on kernel):

- linear kernel: **parametric** method

fixed number of parameters

$$\mathbf{w} \in \mathbb{R}^D$$

- non-linear kernel: **non-parametric** method

number of parameters  $\alpha_i$  increases with number of data points

$$\alpha \in \mathbb{R}^N$$

# Spectral Clustering

Slides by Ulrike von Luxburg (MPI biol. Kybernetik, Tübingen)

\* defered to next week \*

## Further Reading

- Matrix calculus: <http://www.cs.toronto.edu/~roweis/notes.html>
- Multivariate normal distribution:  
[http://en.wikipedia.org/wiki/Multivariate\\_normal\\_distribution](http://en.wikipedia.org/wiki/Multivariate_normal_distribution)
- Shrinkage: **Inadmissibility of the Usual Estimator for the Mean of a Multivariate Normal Distribution.** *Charles Stein.* Proc. Third Berkeley Symp. on Math. Statist. and Prob., Vol. 1 (Univ. of Calif. Press, 1956), 197-206. <http://projecteuclid.org/euclid.bsmmsp/1200501656>
- Least Squares and Logistic Regression: **The Elements of Statistical Learning.** *Hastie, Tibshirani and Friedman.* Springer-Verlag, 2001.
- GPs: <http://www.gaussianprocess.org/>
- SVMs: <http://www.svms.org/tutorials/>
- Kernels and Kernel Machines:
  - **Learning with Kernels.** *Bernhard Schölkopf and Alex Smola.* MIT Press, Cambridge, MA, 2002.
  - <http://www.kernel-machines.org/>
- Spectral Clustering: **A Tutorial on Spectral Clustering.** *Ulrike von Luxburg.* Statistics and Computing 17(4): 395-416, 2007.
- Any statistical term: <http://en.wikipedia.org/>  
(eg Shrinkage estimation of covariance matrices:  
[http://en.wikipedia.org/wiki/Estimation\\_of\\_covariance\\_matrices](http://en.wikipedia.org/wiki/Estimation_of_covariance_matrices))

## Schedule "Machine Learning II" SS'08

- 22.04. Semi-Supervised Learning
- 29.04. Kernels for Structured Data
- 06.05. Applications in Intrusion Detection
- 13.05. Text Mining
- 20.05. Bioinformatics
- 27.05. Optimization for SVMs and Math Programs
- 03.06. Large Scale Optimization
- 10.06. Relevant Dimensionality Estimation
- 17.06. Boosting and Ensemble Methods
- 24.06. Boosting and SVMs
- 01.07. Hidden Markov Models
- 08.07. Structured Output SVMs, Conditional Random Fields
- 15.07. Graphical Models