

Machine Intelligence 1

2.2 Support Vector Machines

Prof. Dr. Klaus Obermayer

Fachgebiet Neuronale Informationsverarbeitung (NI)

WS 2016/2017

2.2.1 Structural Risk Minimization

A bound on the generalization error

- finite samples: bound on the generalization error
(c.f. Statistical Learning Theory, result 3)

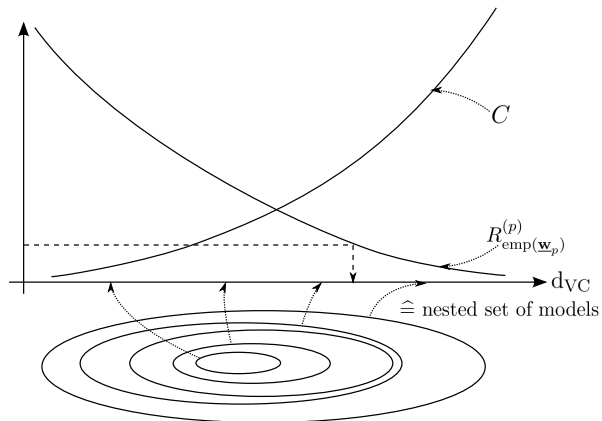
$$P\left\{\sup_{\underline{\mathbf{w}} \in \Lambda} \left| R_{(\underline{\mathbf{w}})} - R_{\text{emp}(\underline{\mathbf{w}})}^{(p)} \right| > \eta \right\} < \underbrace{4 \exp\left(G_{(2p)}^{\Lambda} - p\left(\eta - \frac{1}{p}\right)^2\right)}_{\stackrel{!}{=} \epsilon}$$

- with probability larger than $1 - \epsilon$ we obtain:

$$R_{(\underline{\mathbf{w}})} < \underbrace{R_{\text{emp}(\underline{\mathbf{w}})}^{(p)}}_{\text{empirical error}} + \underbrace{\left(\frac{G_{(2p)}^{\Lambda} - \ln \frac{\epsilon}{4}}{p}\right)^{\frac{1}{2}} + \frac{1}{p}}_{\text{complexity term } C}$$

- For a given ϵ , the complexity term C only depends on p and d_{VC} .

A bound on the generalization error



underfitting \leftarrow ... appropriate model complexity ... \rightarrow overfitting

Structural Risk Minimization (SRM)

$$R_{(\underline{\mathbf{w}})} < R_{\text{emp}(\underline{\mathbf{w}})}^{(p)} + C(p, d_{\text{VC}})$$

- Minimize complexity $C(p, d_{\text{VC}})$ of the model class while keeping the empirical error $R_{\text{emp}(\underline{\mathbf{w}})}^{(p)}$ bounded.
- SRM-learning is consistent (*cf. Vapnik 1998, chapter 6.3*)

2.2.2 Perceptrons Revisited

Canonical hyperplanes

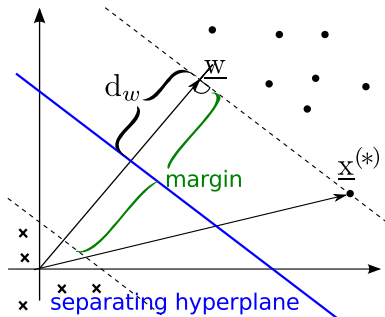
- **data representation** binary classification: $\underline{\mathbf{x}} \in \mathbb{R}^N$, $y_T \in \{-1, +1\}$
- **model class**: connectionist neurons $y = \text{sign}(\underline{\mathbf{w}}^T \underline{\mathbf{x}} + b)$
- parameters of the separating hyperplane $\underline{\mathbf{w}}^T \underline{\mathbf{x}} + b = 0$ are not unique

- data dependent normalization

$$\min_{\alpha=1,\dots,p} \left| \underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right| \stackrel{!}{=} 1$$

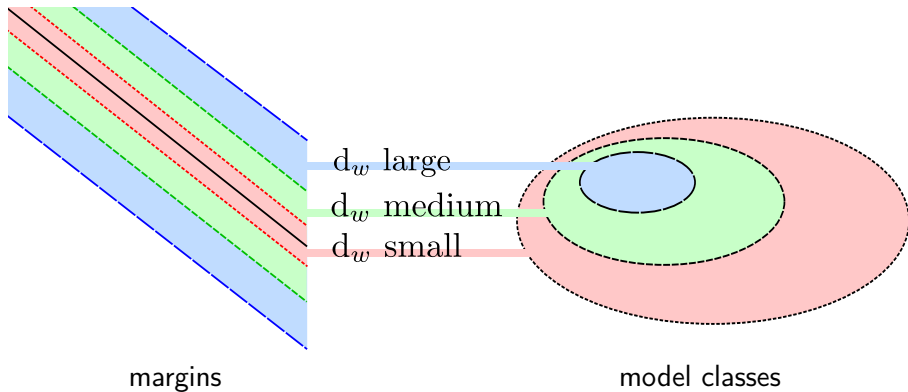
- norm. distance to closest point $\underline{\mathbf{x}}^{(*)}$

$$d_w = \frac{1}{\|\underline{\mathbf{w}}\|} \left| \underline{\mathbf{w}}^T \underline{\mathbf{x}}^* + b \right| \leq \frac{1}{\|\underline{\mathbf{w}}\|}$$

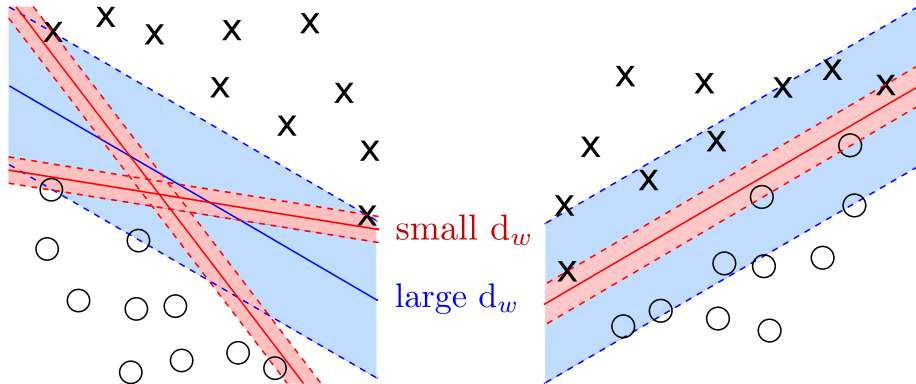


- The minimum normalized distance to the hyperplane is called **margin**.

Nested set of models



Margins and the capacity of the model class



■ larger (minimal) margin \Rightarrow smaller model capacity

Margins and the VC dimension

Theorem (Vapnik, 1998)

$$d_{VC} \leq \min \left(\left\lceil \frac{d_R^2}{d_w^2} \right\rceil, N \right) + 1$$

N : dimension of feature space

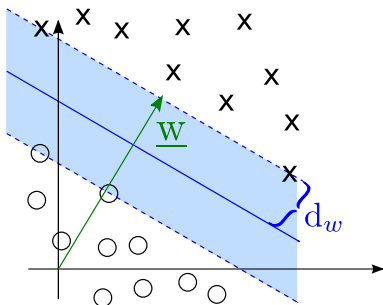
d_w : lower bound of the margin

d_r : range of \underline{x} , $\underline{x} \leq d_R$, for $P(\underline{x}) \neq 0$

- $\frac{d_R^2}{d_w^2}$ is independent of the dimension N of feature space

2.2.3 Learning by Structural Risk Minimization

The primal optimization problem



$$y(\underline{\mathbf{x}}; \underline{\mathbf{w}}) = \text{sign}(\underline{\mathbf{w}}^\top \underline{\mathbf{x}} + b)$$

$$d_w = \frac{1}{\|\underline{\mathbf{w}}\|} \stackrel{!}{=} \max$$

$$\frac{1}{2} \|\underline{\mathbf{w}}\|^2 \stackrel{!}{=} \min$$

(minimize the capacity...)

$$\text{s.t. } y_T^{(\alpha)} \left(\underline{\mathbf{w}}^\top \underline{\mathbf{x}}^{(\alpha)} + b \right) \geq 1, \quad \forall \alpha, \quad (\dots \text{ for zero training error})$$

- inequalities ensure normalization of the weight vector $\underline{\mathbf{w}}$

The method of Lagrange multipliers

$$\underbrace{f_0(\underline{\mathbf{x}}) \stackrel{!}{=} \min}_{\text{minimization}} \quad \text{and} \quad \underbrace{f_k(\underline{\mathbf{x}}) \leq 0, \quad k = 1, \dots, m}_{\text{constraints}}$$

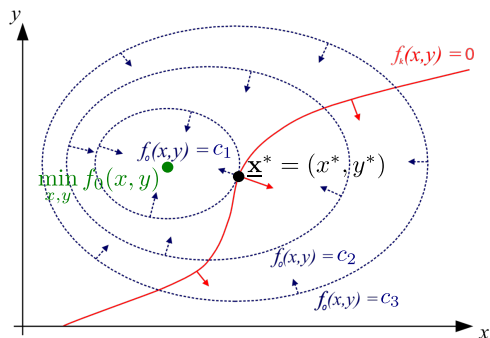
$$L(\underline{\mathbf{x}}, \{\lambda_k\}) \stackrel{!}{=} f_0(\underline{\mathbf{x}}) + \sum_{k=1}^m \lambda_k f_k(\underline{\mathbf{x}}), \quad \lambda_k \geq 0, \quad \forall k \in \{1, \dots, m\}$$

Theorem (Kuhn and Tucker)

Let $A \subset \mathbb{R}^N$ be a convex subset and f_k be convex functions. If there *exists* at least one solution $\underline{\mathbf{x}} \in A$ that satisfies all constraints $f_k(\underline{\mathbf{x}}) \leq 0, \forall k$, then the solution $\underline{\mathbf{x}}^*$ of the constrained optimization problem is given by the saddle point of the Lagrangian, i.e.

$$\min_{\underline{\mathbf{x}} \in A} L(\underline{\mathbf{x}}, \{\lambda_k^*\}) = L(\underline{\mathbf{x}}^*, \{\lambda_k^*\}) = \max_{\lambda_k \geq 0} L(\underline{\mathbf{x}}^*, \{\lambda_k\})$$

The values of the Lagrange multipliers



$$L(\underline{\mathbf{x}}, \{\lambda_k\}) \stackrel{!}{=} f_0(\underline{\mathbf{x}}) + \sum_{k=1}^m \lambda_k f_k(\underline{\mathbf{x}}), \quad \lambda_k \geq 0, \quad \forall k \in \{1, \dots, m\}$$

■ at minimum $\underline{\mathbf{x}}^*$ of boundary $f_k(\underline{\mathbf{x}}) = 0$: $\left. \frac{\partial f_0}{\partial \underline{\mathbf{x}}} \right|_{\underline{\mathbf{x}}^*} \propto - \left. \frac{\partial f_k}{\partial \underline{\mathbf{x}}} \right|_{\underline{\mathbf{x}}^*}$

- $f_k(\underline{\mathbf{x}}^*) = 0 \Rightarrow \lambda_k > 0$ (solution **on** boundary)
- $f_k(\underline{\mathbf{x}}^*) < 0 \Rightarrow \lambda_k = 0$ (solution **behind** boundary)

Application to the primal problem of SRM

- binary classification with linear connectionist neuron

$$f_0(\underline{\mathbf{w}}, b) = \frac{1}{2} \|\underline{\mathbf{w}}\|^2$$

$$f_\alpha(\underline{\mathbf{w}}, b) = -\left\{ y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) - 1 \right\} \leq 0, \quad \forall \alpha \in \{1, \dots, p\}$$

Lagrangian

$$L(\underline{\mathbf{w}}, b, \{\lambda_\alpha\}) = \frac{1}{2} \|\underline{\mathbf{w}}\|^2 - \sum_{\alpha=1}^p \lambda_\alpha \left\{ y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) - 1 \right\}$$

$$\min_{\underline{\mathbf{w}}, b} L(\underline{\mathbf{w}}, b, \{\lambda_\alpha^*\}) = L(\underline{\mathbf{w}}^*, b^*, \{\lambda_\alpha^*\}) = \max_{\lambda_\alpha \geq 0} L(\underline{\mathbf{w}}^*, b^*, \{\lambda_\alpha\})$$

$\underline{\mathbf{w}}, b$: “primal” variables

λ_α : “dual” variables

(solution see blackboard)

The dual problem

$$\underline{\mathbf{w}}^* = \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)}$$

$$L = -\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_{\alpha} \lambda_{\beta} y_T^{(\alpha)} y_T^{(\beta)} \underbrace{\left(\underline{\mathbf{x}}^{(\alpha)} \right)^{\top} \underline{\mathbf{x}}^{(\beta)}}_{\circledast} + \sum_{\alpha=1}^p \lambda_{\alpha} \stackrel{!}{=} \max_{\{\lambda_{\alpha}\}}$$

$$\lambda_{\alpha} \geq 0, \quad \forall \alpha \in \{1, \dots, p\}, \quad \text{and} \quad \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} = 0 \quad (\text{constraints})$$

■ solved numerically using “sequential minimal optimization” (SMO)

The optimal classifier

- connectionist neuron classifier

$$y(\underline{\mathbf{x}}) = \text{sign}(\underline{\mathbf{w}}^\top \underline{\mathbf{x}} + b)$$

- When $\{\lambda_\alpha^*\}_{\alpha=1}^p$ are known, we can compute

$$\underline{\mathbf{w}}^* = \sum_{\alpha=1}^p \lambda_\alpha^* y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)},$$

- and the classifier is thus

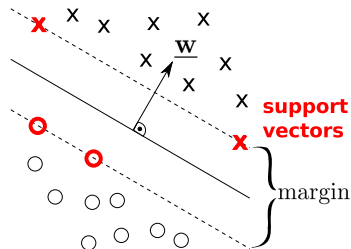
$$y(\underline{\mathbf{x}}) = \text{sign}\left(\sum_{\alpha=1}^p \lambda_\alpha^* y_T^{(\alpha)} \underbrace{(\underline{\mathbf{x}}^{(\alpha)})^\top \underline{\mathbf{x}}}_{(*)} + b^*\right).$$

Support Vectors

- only constraints $f = 0$ on the boundary have $\lambda_\alpha \neq 0$:

$$f_\alpha = -\left\{y_T^{(\alpha)}\left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b\right) - 1\right\} \stackrel{!}{=} 0$$

- constraint $f_\alpha = 0$ implies a normalized distance $d_\alpha = 1$
- these **support vectors** $\underline{\mathbf{x}}^{(\alpha)}$ are thus **on the margin**



Bias calculation

- for all support vectors $\underline{\mathbf{x}}^\alpha \in SV$ **on** the margin holds

$$b^* = y_T^{(\alpha)} - \underline{\mathbf{w}}^{*\top} \underline{\mathbf{x}}^{(\alpha)}.$$

- compute bias b^* as average over all $\underline{\mathbf{x}}^{(\alpha)} \in SV$

$$b^* = \frac{1}{\#SV} \sum_{\alpha \in SV} \left(y_T^{(\alpha)} - \sum_{\beta \in SV} \lambda_\beta y_T^{(\beta)} \underbrace{\left(\underline{\mathbf{x}}^{(\beta)} \right)^T \underline{\mathbf{x}}^{(\alpha)}}_{\circledast} \right)$$

Support Vector Machines (SVM)

- **perceptrons** $\hat{y}(\underline{\mathbf{x}}) = \text{sign}(\underline{\mathbf{w}}^T \underline{\mathbf{x}} + b)$ trained by SRM are called SVM
- weights and threshold are calculated by solving the **dual optimization problem** for Lagrange multipliers $\{\lambda_\alpha \geq 0\}_{\alpha=1}^p$

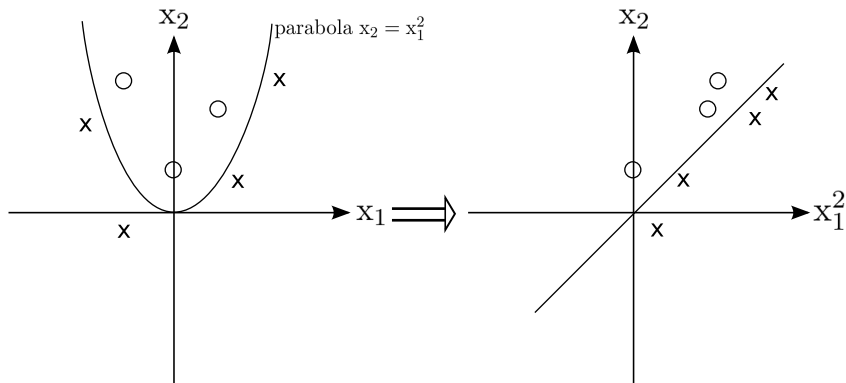
$$\max L(\{\lambda_\alpha\}) = -\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_\alpha \lambda_\beta y_T^{(\alpha)} y_T^{(\beta)} \underbrace{(\underline{\mathbf{x}}^{(\alpha)})^T \underline{\mathbf{x}}^{(\beta)}}_{(*)} + \sum_{\alpha=1}^p \lambda_\alpha$$

- **SVM solution is** $\hat{y}(\underline{\mathbf{x}}) = \text{sign} \left(\sum_{\alpha \in \text{SV}} \lambda_\alpha y_T^{(\alpha)} \underbrace{(\underline{\mathbf{x}}^{(\alpha)})^T \underline{\mathbf{x}}}_{(*)} + b \right)$

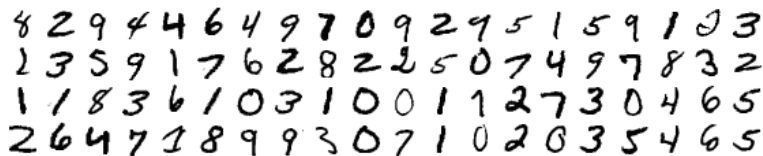
$$\text{with } b = \frac{1}{\#_{\text{SV}}} \sum_{\alpha \in \text{SV}} \left(y_T^{(\alpha)} - \sum_{\beta \in \text{SV}} \lambda_\beta y_T^{(\beta)} \underbrace{(\underline{\mathbf{x}}^{(\beta)})^T \underline{\mathbf{x}}^{(\alpha)}}_{(*)} \right)$$

2.2.4 SRM Learning for Non-linear Classification Boundaries

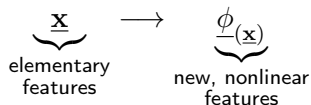
Transformation of feature space



Transformation of feature space



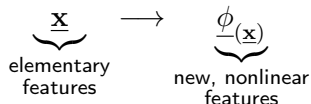
- feature space: monomials of degree n



- $n = 10$ and N pixel values $x_i \Rightarrow N^{10}$ monomials
- SVM requires only **scalar products** $\phi(\underline{\mathbf{x}})^\top \phi(\underline{\mathbf{x}'})$ in feature space

The kernel trick

- “project” data implicitly in high dimensional feature space $\underline{\phi}$



- SVM requires only scalar products $\underline{\phi(\mathbf{x})}^T \underline{\phi(\mathbf{x})}$ in feature space
- replace scalar products with **kernel function** $\underline{\phi(\mathbf{x})}^T \underline{\phi(\mathbf{x}')} \rightarrow K(\underline{\mathbf{x}}, \underline{\mathbf{x}}')$

Mercer's theorem

- let \mathcal{X} be a *compact* subset of \mathbb{R}^N
- let $K : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}, K \in L_\infty$, be a symmetric function ("kernel")
- let $T_K : L_2(\mathcal{X}) \rightarrow L_2(\mathcal{X})$ be the linear convolution operator

$$T_K[f]_{(\underline{\mathbf{x}})} := \int_{\mathcal{X}} K_{(\underline{\mathbf{x}}, \underline{\mathbf{x}'})} f_{(\underline{\mathbf{x}'})} d\underline{\mathbf{x}'}$$

- let $\lambda_i \in \mathbb{R}$ be *eigenvalues* and $\psi_{i(\underline{\mathbf{x}})} \in L_2(\mathcal{X})$ *eigenfunctions* of T_K

Mercer's theorem

Every **positive semi-definite** kernel K corresponds to a **scalar product** $K(\underline{\mathbf{x}}, \underline{\mathbf{x}'}) = \underline{\phi}_{(\underline{\mathbf{x}})}^\top \underline{\phi}_{(\underline{\mathbf{x}'})}$ in the feature space spanned by $\phi_{i(\underline{\mathbf{x}})} = \sqrt{\lambda_i} \psi_{i(\underline{\mathbf{x}})}$.

Kernel properties

symmetric kernels

orthonormal eigenfunctions:

$$\int_{\mathcal{X}} \psi_{i(\underline{\mathbf{x}})} \psi_{j(\underline{\mathbf{x}})} d\underline{\mathbf{x}} = \delta_{ij}, \forall i, j \in \mathbb{N}$$

positive semi-definite kernels

all eigenvalues λ_i are **non-negative**:

$$\iint_{\mathcal{X} \times \mathcal{X}} K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') f(\underline{\mathbf{x}}) f(\underline{\mathbf{x}}') d\underline{\mathbf{x}} d\underline{\mathbf{x}}' \geq 0, \quad \forall f \in L_2(\mathcal{X}), \quad (\text{positive semi-definite})$$

The feature space induced by kernels may be infinite dimensional!

Typical kernel functions

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = (\underline{\mathbf{x}}^T \underline{\mathbf{x}}' + 1)^d$$

polynomial kernel of degree d
→ image processing: pixel correlations

Typical kernel functions

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = (\underline{\mathbf{x}}^T \underline{\mathbf{x}}' + 1)^d$$

polynomial kernel of degree d
→ image processing: pixel correlations

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \exp \left\{ - \frac{(\underline{\mathbf{x}} - \underline{\mathbf{x}}')^2}{2\sigma^2} \right\}$$

RBF-kernel with range σ
→ infinite dimensional feature space

Typical kernel functions

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = (\underline{\mathbf{x}}^T \underline{\mathbf{x}}' + 1)^d$$

polynomial kernel of degree d
→ image processing: pixel correlations

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \exp \left\{ - \frac{(\underline{\mathbf{x}} - \underline{\mathbf{x}}')^2}{2\sigma^2} \right\}$$

RBF-kernel with range σ
→ infinite dimensional feature space

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \tanh \{ \kappa \underline{\mathbf{x}}^T \underline{\mathbf{x}}' + \theta \}$$

neural network kernel with parameters κ and θ → not positive definite!

Typical kernel functions

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = (\underline{\mathbf{x}}^T \underline{\mathbf{x}}' + 1)^d$$

polynomial kernel of degree d
→ image processing: pixel correlations

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \exp \left\{ - \frac{(\underline{\mathbf{x}} - \underline{\mathbf{x}}')^2}{2\sigma^2} \right\}$$

RBF-kernel with range σ
→ infinite dimensional feature space

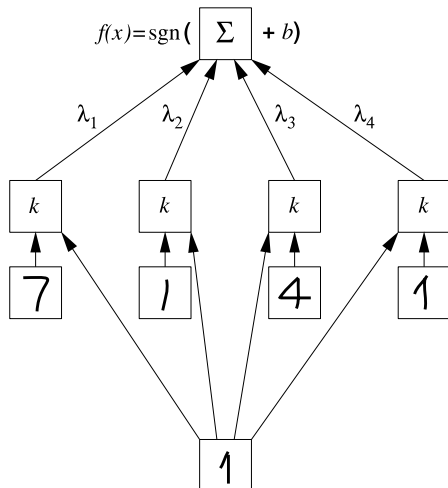
$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \tanh \{ \kappa \underline{\mathbf{x}}^T \underline{\mathbf{x}}' + \theta \}$$

neural network kernel with parameters κ and θ → not positive definite!

$$K(\underline{\mathbf{x}}, \underline{\mathbf{x}}') = \frac{1}{(\|\underline{\mathbf{x}} - \underline{\mathbf{x}}'\|^2 + \epsilon^2)^{N/2}}$$

Plummer kernel with parameter ϵ
→ scale invariant kernel

SVM with kernels



classification

$$f(x) = \text{sgn} \left(\sum_i \lambda_i k(x, x_i) + b \right)$$

weights

comparison: e.g. $k(x, x_i) = (x \cdot x_i)^d$ support vectors $x_1 \dots x_4$ $k(x, x_i) = \exp(-\|x - x_i\|^2 / c)$ $k(x, x_i) = \tanh(\kappa(x \cdot x_i) + \theta)$ input vector x

see Schölkopf & Smola (2001, p. 202)

Comments

- Mercer's theorem can be used to "kernelize" many different linear methods, both supervised or unsupervised.
 - Fisher discriminant analysis
 - principal component analysis (*see MI 2*)
 - k -means clustering & self-organizing maps
 - canonical correlation analysis

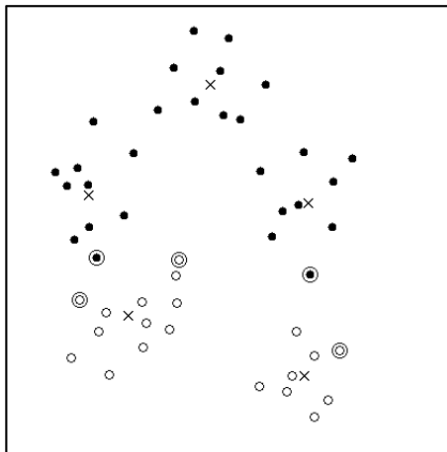
Comments

- SVM vs. RBF networks
- RBF network for classification
 - 5 Gaussian bases (\times)

$$y(\underline{\mathbf{x}}) = \text{sign} \left(\sum_{i=1}^5 w_i \exp\left(\frac{1}{2\sigma_i^2} \|\underline{\mathbf{x}} - \underline{\mathbf{t}}_i\|^2\right) \right)$$

- SVM with Gaussian kernel
 - 5 support vectors $\underline{\mathbf{x}}_i$ (\circ)
 - $a_i = \lambda_\alpha y_T^{(\alpha)}$, for $\underline{\mathbf{x}}_i = \underline{\mathbf{x}}^{(\alpha)}$

$$y(\underline{\mathbf{x}}) = \text{sign} \left(\sum_{i=1}^5 a_i \exp\left(\frac{1}{2\sigma^2} \|\underline{\mathbf{x}} - \underline{\mathbf{x}}_i\|^2\right) + b \right)$$

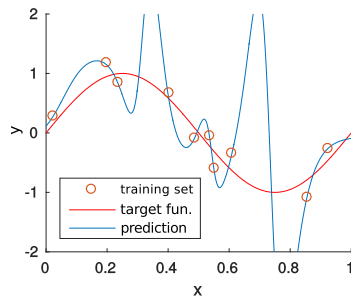


see Schölkopf et al. (1997), Schölkopf & Smola (2001, p. 204)

2.2.5 The C-Support Vector Machine

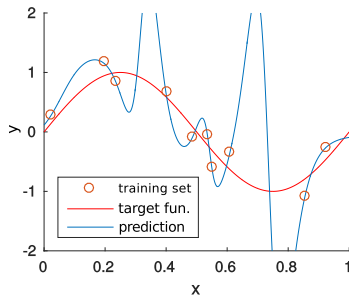
Classification of non-separable problems

- real-world problems are typically non-separable
- incomplete feature sets & noise
- perfect separation of the training set \leadsto overfitting



Classification of non-separable problems

- real-world problems are typically non-separable
- incomplete feature sets & noise
- perfect separation of the training set \leadsto overfitting



consequences

$$R(\underline{\mathbf{w}}) \leq R_{\text{emp}}^{(p)}(\underline{\mathbf{w}}) + C(p, d_{\text{VC}})$$

- finite training error $R_{\text{emp}}^{(p)} \neq 0$
- trade-off between minimization of the training error and the capacity of the model class

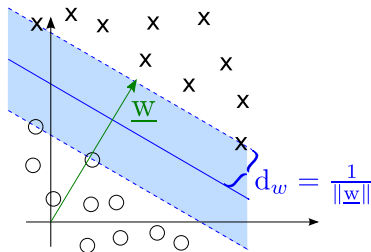
The primal problem

$$\frac{1}{2} \|\underline{\mathbf{w}}\|^2 \quad \stackrel{!}{=} \min \quad \left\{ \begin{array}{l} \text{minimize upper bound on VC dimension} \end{array} \right.$$

constraints ($\forall \alpha$):

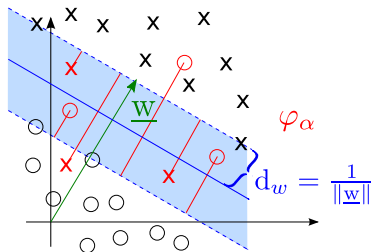
$$y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) \geq 1$$

normalization & correct classification
of all data points



constraints ($\forall \alpha$): (C : regularization parameter)

$$\begin{aligned} y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) &\geq 1 - \varphi_\alpha && \text{normalization \& correct classification} \\ &&& \text{of all data points for } \varphi_\alpha=0 \\ \varphi_\alpha &\geq 0 && \text{"margin errors" for } \varphi_\alpha \neq 0 \end{aligned}$$



Dual problem of the C-SVM

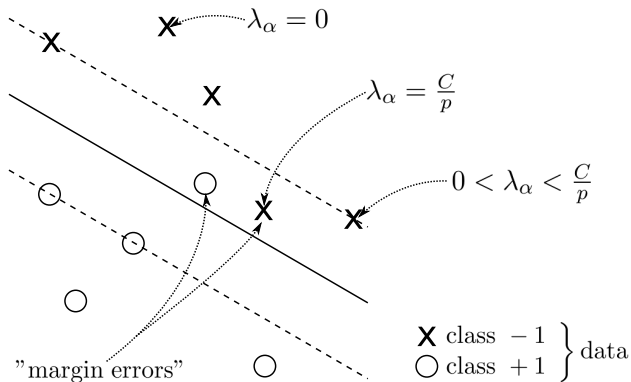
■ Objective

$$-\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_{\alpha} \lambda_{\beta} y_T^{(\alpha)} y_T^{(\beta)} \underbrace{\left(\underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\mathbf{x}}^{(\beta)}}_{\text{kernel function}} + \sum_{\alpha=1}^p \lambda_{\alpha} \stackrel{!}{=} \max_{\{\lambda_{\alpha}\}_{\alpha=1}^p}$$

■ Constraints:

$$\sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} = 0 \qquad 0 \leq \underbrace{\lambda_{\alpha}}_{\substack{\text{difference to} \\ \text{separable case}}} \leq \frac{C}{p}$$

Margin and support vectors



The C-SVM classifier

$$\underline{\mathbf{w}} = \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)} \quad \leadsto \lambda_{\alpha} \neq 0 \text{ only for support vectors } SV$$

The C-SVM classifier

$$\underline{\mathbf{w}} = \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)} \quad \leadsto \lambda_{\alpha} \neq 0 \text{ only for support vectors } SV$$

$$b = \frac{1}{\#SV_{<}} \sum_{\alpha \in SV_{<}} \left(y_T^{(\alpha)} - \sum_{\beta \in SV} \lambda_{\beta} y_T^{(\beta)} \underbrace{\left(\underline{\mathbf{x}}^{(\beta)} \right)^T \underline{\mathbf{x}}^{(\alpha)}}_{\text{kernel!}} \right)$$

$SV_{<}$: SV s with $\lambda_{\alpha} < \frac{C}{p}$ (SV s on the margin)

The C-SVM classifier

$$\underline{\mathbf{w}} = \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)} \quad \leadsto \lambda_{\alpha} \neq 0 \text{ only for support vectors } SV$$

$$b = \frac{1}{\#SV_{<}} \sum_{\alpha \in SV_{<}} \left(y_T^{(\alpha)} - \sum_{\beta \in SV} \lambda_{\beta} y_T^{(\beta)} \underbrace{\left(\underline{\mathbf{x}}^{(\beta)} \right)^T \underline{\mathbf{x}}^{(\alpha)}}_{\text{kernel!}} \right)$$

$SV_{<}$: SV s with $\lambda_{\alpha} < \frac{C}{p}$ (SV s on the margin)

Classifier

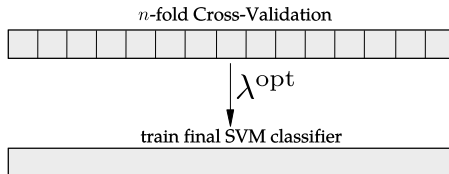
$$\hat{y}(\underline{\mathbf{x}}) = \text{sign}(\underline{\mathbf{w}}^T \underline{\mathbf{x}} + b) = \text{sign} \left(\sum_{\alpha \in SV} \lambda_{\alpha} y_T^{(\alpha)} \underbrace{\left(\underline{\mathbf{x}}^{(\alpha)} \right)^T \underline{\mathbf{x}}}_{\text{kernel function}} + b \right)$$

Validation & selection of hyperparameters

- validation and model selection w.r.t. **0-1 loss**

$$e(\underline{\mathbf{x}}^{(\alpha)}, y_T^{(\alpha)}) = \begin{cases} 0 & , \text{if } \hat{y}(\underline{\mathbf{x}}^{(\alpha)}) = y_T^{(\alpha)} \\ 1 & , \text{otherwise} \end{cases}$$

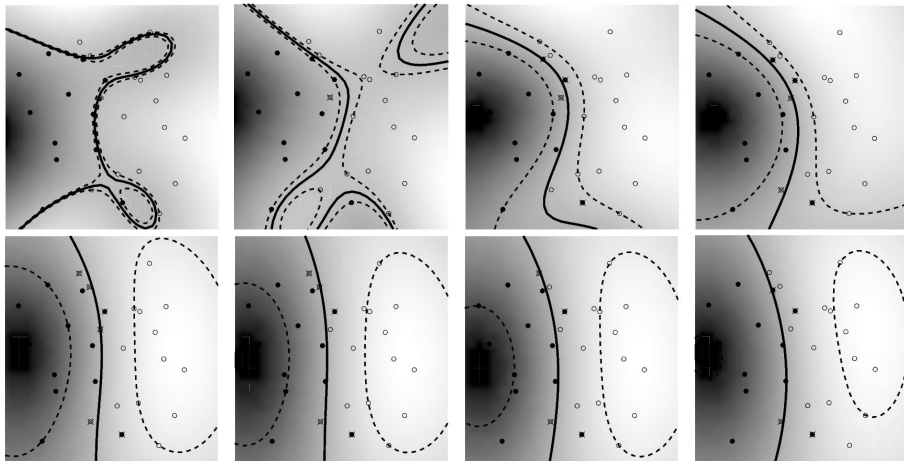
- hyper-parameter selection (C, σ, \dots) by n -fold **cross-validation**



- validation on hold-out **validation set**



SVM and overfitting



(related ν -SVM with $\nu \in \{0.1, 0.2, \dots, 0.8\}$ and RBF kernel)

see Schölkopf & Smola (2001, p. 207)

2.2.6 Sequential Minimal Optimization

The dual problem

$$\begin{aligned}
 & -\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_{\alpha} \lambda_{\beta} y_T^{(\alpha)} y_T^{(\beta)} K_{(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})} + \sum_{\alpha=1}^p \lambda_{\alpha} \stackrel{!}{=} \max_{\{\lambda_{\alpha}\}_{\alpha=1}^p} \\
 & \text{s.t.} \quad \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} = 0, \quad 0 \leq \lambda_{\alpha} \leq \frac{C}{p}.
 \end{aligned}$$

The dual problem

$$-\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_{\alpha} \lambda_{\beta} y_T^{(\alpha)} y_T^{(\beta)} K_{\alpha\beta} + \sum_{\alpha=1}^p \lambda_{\alpha} \stackrel{!}{=} \max_{\{\lambda_{\alpha}\}_{\alpha=1}^p}$$

$$\text{s.t.} \quad \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} = 0, \quad 0 \leq \lambda_{\alpha} \leq \frac{C}{p}.$$

The Gram matrix **K**

$$K_{\alpha\beta} = K_{(\underline{\mathbf{x}}^{(\alpha)}, \underline{\mathbf{x}}^{(\beta)})}$$

	1	2	3	...	j
1	K_{11}	K_{12}	K_{1j}
2	\vdots	\vdots	K_{23}	...	K_{2j}
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
i	K_{i1}	K_{i2}	K_{ij}

The dual problem

$$-\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_{\alpha} \lambda_{\beta} y_T^{(\alpha)} y_T^{(\beta)} K_{\alpha\beta} + \sum_{\alpha=1}^p \lambda_{\alpha} \stackrel{!}{=} \max_{\{\lambda_{\alpha}\}_{\alpha=1}^p}$$

$$\text{s.t.} \quad \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} = 0, \quad 0 \leq \lambda_{\alpha} \leq \frac{C}{p}.$$

- SVMs operate on pairwise (similarity) data!
- positive definite kernel \rightarrow positive definite Gram matrix **K**
 \Rightarrow well defined optimization problem
- **K** can be pre-computed to speed up subsequent computations.

The SMO procedure

$$-\frac{1}{2} \sum_{\alpha, \beta=1}^p \lambda_{\alpha} \lambda_{\beta} y_T^{(\alpha)} y_T^{(\beta)} K_{\alpha\beta} + \sum_{\alpha=1}^p \lambda_{\alpha} \stackrel{!}{=} \max_{\{\lambda_{\alpha}\}_{\alpha=1}^p}$$

$$\text{s.t.} \quad \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} = 0, \quad 0 \leq \lambda_{\alpha} \leq \frac{C}{p}.$$

while *not converged* **do**

 Choose two Lagrange multipliers $\lambda_{\gamma}, \lambda_{\delta}$.

 Optimize the constrained Lagrangian while changing only λ_{γ} and λ_{δ} .

end

Choosing λ_γ and λ_δ based on KKT

Karush-Kuhn-Tucker conditions (KKT conditions)

$$\underbrace{\left[y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) - 1 + \varphi_\alpha \right]}_{\substack{\text{constraint of the primal problem:} \\ = 0 \text{ for all data points on and} \\ \text{within the margin}}} \underbrace{\lambda_\alpha}_{\substack{\text{Lagrange mul.} \\ = 0 \text{ for all data} \\ \text{points outside} \\ \text{the margin}}} = 0 \quad (\text{KKT})$$

- ① loop over all λ_γ violating KKT-conditions
(and additional "threshold"-conditions due to errors in b)
pick λ_γ for which $\text{KKT} \neq 0$
- ② for this λ_γ : select λ_δ yielding a "large step" towards optimum (general heuristics, difference in relative errors $f(x^{(\alpha)}) - y^{(\alpha)}$ vs. $f(x^{(\beta)}) - y^{(\beta)}$)

Reduced optimization problem

$$\min_{(\lambda_\alpha)} \stackrel{!}{=} \frac{1}{2} \sum_{\alpha\beta} \lambda_\alpha \lambda_\beta y_T^{(\alpha)} y_T^{(\beta)} K_{\alpha\beta} - \sum_{\alpha} \lambda_\alpha$$

Reduced optimization problem

$$\begin{aligned}
 \min_{(\lambda_\alpha)} & \stackrel{!}{=} \frac{1}{2} \sum_{\alpha\beta} \lambda_\alpha \lambda_\beta y_T^{(\alpha)} y_T^{(\beta)} K_{\alpha\beta} - \sum_{\alpha} \lambda_\alpha \\
 \min_{(\lambda_\delta, \lambda_\gamma)} & \stackrel{!}{=} \frac{1}{2} \left[\underbrace{\lambda_\gamma^2 \left(y_T^{(\gamma)} \right)^2 K_{\gamma\gamma}}_{=1} + \underbrace{\lambda_\delta^2 \left(y_T^{(\delta)} \right)^2 K_{\delta\delta}}_{=1} + 2\lambda_\gamma \lambda_\delta \underbrace{y_T^{(\gamma)} y_T^{(\delta)} K_{\gamma\delta}}_{Q_{\gamma\delta}} \right] \\
 & + \lambda_\gamma \left[\underbrace{\sum_{\beta \neq \delta, \gamma} \lambda_\beta y_T^{(\gamma)} y_T^{(\beta)} K_{\gamma\beta}}_{C_\gamma} - 1 \right] + \lambda_\delta \left[\underbrace{\sum_{\beta \neq \gamma, \delta} \lambda_\beta y_T^{(\delta)} y_T^{(\beta)} K_{\delta\beta}}_{C_\delta} - 1 \right] + \text{const}_{(\lambda_\delta, \lambda_\gamma)} \\
 \min_{(\lambda_\delta, \lambda_\gamma)} & \stackrel{!}{=} \frac{1}{2} \left[\lambda_\gamma^2 Q_{\gamma\gamma} + \lambda_\delta^2 Q_{\delta\delta} + 2\lambda_\gamma \lambda_\delta Q_{\gamma\delta} \right] + C_\gamma \lambda_\gamma + C_\delta \lambda_\delta
 \end{aligned}$$

Sequential Minimal Optimization (SMO)

- optimize

$$\min_{(\lambda_\delta, \lambda_\gamma)} \stackrel{!}{=} \frac{1}{2} \left[\lambda_\gamma^2 Q_{\gamma\gamma} + \lambda_\delta^2 Q_{\delta\delta} + 2\lambda_\gamma \lambda_\delta Q_{\gamma\delta} \right] + C_\gamma \lambda_\gamma + C_\delta \lambda_\delta$$

- under the following "box" and "equality" constraints

$$0 \leq \lambda_{\gamma,\delta} \leq \frac{C}{p}, \quad (\text{i})$$

$$\underbrace{\lambda_\gamma + \frac{y_T^{(\delta)}}{y_T^{(\gamma)}} \lambda_\delta}_s = - \underbrace{\frac{1}{y_T^{(\gamma)}} \sum_{\beta \neq \gamma, \delta} \lambda_\beta y_T^{(\beta)}}_d \Rightarrow \lambda_\gamma + s \lambda_\delta = -d \quad (\text{ii})$$

- Analytical solution: Schoelkopf & Smola, p. 308
- Pseudocode: Schoelkopf & Smola, p. 313
- Software: www.csie.ntu.edu.tw/~cjlin/libsvm/
(also covers multiclass problems, support vector regression, one-class SVMs)

Remarks

Sequential Minimal Optimization (SMO) ...

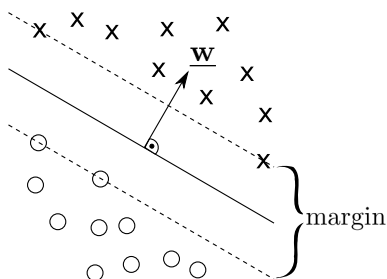
- ...exploits that for 2 constraints the optimization problem can be solved analytically
- ...needs little memory (\approx number of datapoints)
- ...can be much faster than other algorithms
- ...convergence speed depends on rules to select the λ_i
 \leadsto good heuristics are important

End of Section 2.2

the following slides contain

OPTIONAL MATERIAL

Classification margin



- **margin**: minimal (normalized) distance to hyperplane $d^{\min} = \frac{1}{\|\underline{w}\|}$
- large margins have low *ambiguity* \Rightarrow low VC-dimension

The solution of the primal problem

Lagrangian

$$L = \frac{1}{2} \|\underline{\mathbf{w}}\|^2 - \sum_{\alpha=1}^p \lambda_{\alpha} \left\{ y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) - 1 \right\}$$

- setting derivative w.r.t. weights w_l to zero: $\underline{\mathbf{w}} = \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)}$

$$\frac{\partial L}{\partial w_l} = w_l - \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} x_l^{(\alpha)} \stackrel{!}{=} 0$$

The solution of the primal problem

Lagrangian

$$L = \frac{1}{2} \|\underline{\mathbf{w}}\|^2 - \sum_{\alpha=1}^p \lambda_{\alpha} \left\{ y_T^{(\alpha)} \left(\underline{\mathbf{w}}^T \underline{\mathbf{x}}^{(\alpha)} + b \right) - 1 \right\}$$

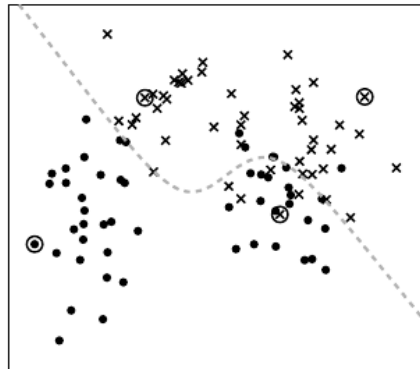
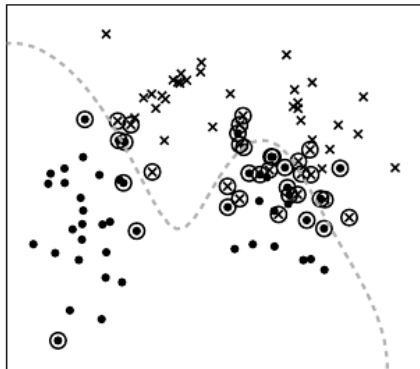
- setting derivative w.r.t. weights w_l to zero: $\underline{\mathbf{w}} = \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \underline{\mathbf{x}}^{(\alpha)}$

$$\frac{\partial L}{\partial w_l} = w_l - \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} x_l^{(\alpha)} \stackrel{!}{=} 0$$

- setting derivative w.r.t. b to zero

$$\frac{\partial L}{\partial b} = - \sum_{\alpha=1}^p \lambda_{\alpha} y_T^{(\alpha)} \stackrel{!}{=} 0$$

Sparse Bayesian Regression: Relevance Vector Machines



see Tipping (2001)