## Introduction to Statistical Machine Learning

### Christfried Webers

Statistical Machine Learning Group NICTA and College of Engineering and Computer Science The Australian National University

> Canberra February – June 2013

(Many figures from C. M. Bishop, "Pattern Recognition and Machine Learning")

#### Introduction to Statistical Machine Learning

Christfried Webers NICTA

The Australian National University



### Overview

Introduction Linear Algebra Probability Linear Regression 1 Linear Regression 2 Linear Classification 1 Linear Classification 2

Neural Networks 1 Neural Networks 2 Kernel Methods Sparse Kernel Methods

Graphical Models 1

Graphical Models 2

Graphical Models 3 Mixture Models and FM 1 Mixture Models and EM 2

Approximate Inference

Sampling

Principal Component Analysis

Sequential Data 1 Sequential Data 2

Combining Models Selected Topics

Discussion and Summary

1of 463

# Part XII

Sparse Kernel Machines

#### Introduction to Statistical Machine Learning

© 2013 Christfried Webers NICTA

The Australian National University



Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class

© 2013
Christfried Webers
NICTA
The Australian National



Sparse Kernel Machines

Maximum Margi Classifiers

Overlapping Cla Distributions

Relevance Vector

- Nonlinear kernels extended our toolbox of methods considerably. Wherever an inner product was used in an algorithms, we can replace it with a kernel.
- Kernels act as a kind of 'similarity' measure and can be defined over graphs, sets, strings, and documents.
- But the kernel matrix is a square matrix with dimensions equal to the number of data points N. In order to calculate it, the kernel function must be evaluated for all pairs of training inputs.
- Sparse Kernel Machines implement learning algorithms which are based on sparse kernels. For prediction, these kernels are only evaluated at a subset of the training data.

Maximum Margin Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

Return to the two-class classification with a linear model

$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}) + b$$

where  $\phi(\mathbf{x})$  is some fixed feature space mapping, and b is the bias.

- Training data are N input vectors  $\mathbf{x}_1, \dots, \mathbf{x}_N$  with corresponding targets  $t_1, \dots, t_N$  where  $t_n \in \{-1, +1\}$ .
- The class of a new point is predicted as  $sign(y(\mathbf{x}))$ .
- Assume there exists a linear decision boundary. That means, there exist w and b such that

$$t_n \operatorname{sign}(y(\mathbf{x}_n)) > 0$$
  $n = 1, \ldots, N.$ 

- The perceptron algorithm can find a solution, but this depends on the initial choice of parameters.
- What is the decision boundary which results in the best generalisation (smallest generalisation error)?



Introduction to Statistical Machine Learning

© 2013 Christfried Webers NICTA The Australian National

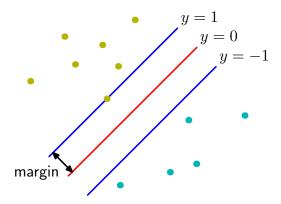


Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class Distributions

 The margin is the smallest distance between the decision boundary and any of the samples. (We will see later why  $y = \pm 1$  on the margin.)



Introduction to Statistical Machine Learning

Christfried Webers NICTA The Australian National



Maximum Margin Classifiers

which maximises the margin.

© 2013 Christfried Webers NICTA

The Australian National University

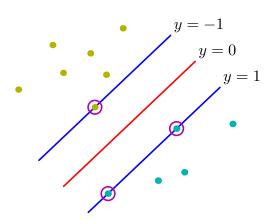


Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview



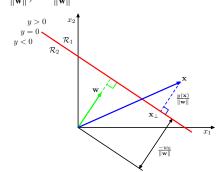
Support Vector Machines choose the decision boundary



Maximum Margin Classifiers

Overlapping Class
Distributions

- w is orthogonal to the decision boundary  $(0 = \mathbf{w}^T \mathbf{x_1} + w_0 = \mathbf{w}^T \mathbf{x_2} + w_0 \text{ and } \mathbf{w}^T (\mathbf{x_1} \mathbf{x_2}) = 0.)$
- Length of the projection of x onto w is  $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|}$ .
- If **x** sits on the decision boundary then  $0 = y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$  and therefore the distance from the origin orthogonal to the decision boundary (and parallel to **w** is  $\frac{\mathbf{w}^T \mathbf{x}}{\|\mathbf{w}\|} = -\frac{w_0}{\|\mathbf{w}\|}$ .
- The distance of x from the decision boundary is therefore  $\frac{\mathbf{w}^T\mathbf{x}}{\|\mathbf{w}\|} (-\frac{w_0}{\|\mathbf{w}\|}) = \frac{y(\mathbf{x})}{\|\mathbf{w}\|}.$



- © 2013 Christfried Webers NICTA The Australian National University
- ISML 2013

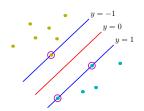
Maximum Margin Classifiers

Relevance Vector

Machines - Overview

- Support Vector Machines choose the decision boundary which maximises the smallest distance to samples in both classes.
- We calculated the distance of a point  $\mathbf{x}$  from the hyperplane  $y(\mathbf{x}) = 0$  as  $|y(\mathbf{x})|/|\mathbf{w}||$ .
- Perfect classification means  $t_n y(\mathbf{x}_n) > 0$  for all n.
- $\bullet$  Thus, the distance of  $\mathbf{x}_n$  from the decision boundary is

$$\frac{t_n y(\mathbf{x}_n)}{\|\mathbf{w}\|} = \frac{t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b)}{\|\mathbf{w}\|}$$





Maximum Margin Classifiers

Overlapping Class Distributions

Machines - Overview

- Support Vector Machines choose the decision boundary which maximises the smallest distance to samples in both classes.
- For the maximum margin solution, solve

$$\arg\max_{\mathbf{w},b} \left\{ \frac{1}{\|\mathbf{w}\|} \min_{n} \left[ t_n \left( \mathbf{w}^T \phi(\mathbf{x}_n) + b \right) \right] \right\}.$$

• How can we solve this?



Maximum Marein Classifiers

Maximum margin solution

$$\underset{\mathbf{w},b}{\operatorname{arg\,max}} \left\{ \frac{1}{\|\mathbf{w}\|} \min_{n} \left[ t_{n} \left( \mathbf{w}^{T} \boldsymbol{\phi}(\mathbf{x}_{n}) + b \right) \right] \right\}$$

- Observation: Rescaling  $\mathbf{w} \to \kappa \mathbf{w}$  and  $b \to \kappa b$  does not change the distance from the hypersurface.
- Scale in such a way that for the closest point

$$t_n\left(\mathbf{w}^T\boldsymbol{\phi}(\mathbf{x}_n)+b\right)=1.$$

 The canonical representation for the decision hyperplane is therefore

$$t_n(\mathbf{w}^T\phi(\mathbf{x}_n)+b)\geq 1 \qquad n=1,\ldots,N.$$

 The constraints are active, if the equality holds, otherwise they are inactive.

Maximum Margin Classifiers

Overlapping Class
Distributions

Relevance Vector Machines - Overview

Maximum margin solution

$$\underset{\mathbf{w},b}{\operatorname{arg\,max}} \left\{ \frac{1}{\|\mathbf{w}\|} \min_{n} \left[ t_{n} \left( \mathbf{w}^{T} \phi(\mathbf{x}_{n}) + b \right) \right] \right\}$$

Transformed once

$$\underset{\mathbf{w},b}{\operatorname{arg\,max}} \left\{ \frac{1}{\|\mathbf{w}\|} \right\} \quad \text{s.t. } t_n \left( \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n) + b \right) \geq 1.$$

Transformed again

$$\arg\min_{\mathbf{w},b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 \right\} \quad \text{s.t. } t_n \left( \mathbf{w}^T \phi(\mathbf{x}_n) + b \right) \ge 1.$$



Maximum Margin Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

 Quadratic Programming (QP) problem: minimise a quadratic function subject to linear constraints.

$$\arg\min_{\mathbf{w},b} \left\{ \frac{1}{2} \|\mathbf{w}\|^2 \right\} \quad \text{s.t. } t_n \left( \mathbf{w}^T \phi(\mathbf{x}_n) + b \right) \ge 1.$$

• Introduce Lagrange multipliers  $\{a_n \ge 0\}$  to get the Lagrangian

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} ||\mathbf{w}||^2 - \sum_{n=1}^{N} a_n \{ t_n (\mathbf{w}^T \phi(\mathbf{x}_n) + b) - 1 \}.$$

NICTA

The Australian National

$$L(\mathbf{w}, b, \mathbf{a}) = \frac{1}{2} ||\mathbf{w}||^2 - \sum_{n=1}^{N} a_n \left\{ t_n \left( \mathbf{w}^T \phi(\mathbf{x}_n) + b \right) - 1 \right\}.$$

• Derivatives with respect to w and b are

$$\mathbf{w} = \sum_{n=1}^{N} a_n t_n \phi(\mathbf{x}_n) \qquad 0 = \sum_{n=1}^{N} a_n t_n$$

Dual representation (again QP problem): Maximise

$$\widetilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\sum_{n=1}^{N} a_n t_n = 0$$

$$a_n \geq 0$$
  $n = 1, \ldots, N.$ 



Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class
Distributions



Maximum Margin Classifiers

Overlapping Clas. Distributions

Relevance Vector

Kernel function

$$k(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}')$$

must imply a positive definite kernel (Gram matrix), otherwise the optimisation problem is not bounded from above.

- Solving the optimisation problem is of the order  $O(N^3)$  in the number of data points N. The original optimisation problem is of order  $O(M^3)$  where M was the number of basis functions.
- Again, kernel allows us to use large or infinite number of basis functions (feature maps), but disadvantegeous for M < N.</li>

Introduction to Statistical

Machine Learning

Maximum Margin Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

$$y(\mathbf{x}) = \sum_{n=1}^{N} a_n t_n k(\mathbf{x}, \mathbf{x}_n) + b$$

The Karush-Kuhn-Tucker (KKT) conditions are

$$a_n \ge 0$$

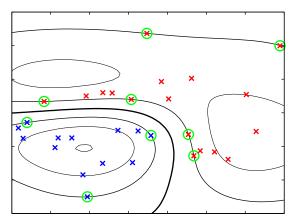
$$t_n y(\mathbf{x}_n) - 1 \ge 0$$

$$a_n \{t_n y(\mathbf{x}_n) - 1\} = 0.$$

- Therefore, either  $a_n = 0$  or  $t_n y(\mathbf{x}_n) 1 = 0$ .
- If  $a_n = 0$ , no contribution to the prediction!
- After training, use only the set S of points for which  $a_n > 0$  (and therefore  $t_n y(\mathbf{x}_n) = 1$ ) holds.
- S contains only the support vectors.

## Maximum Margin Classifiers - Support Vectors

• S contains only the support vectors.



Decision and margin boundaries for a two-class problem using Gaussian kernel functions.

Introduction to Statistical Machine Learning

© 2013 Christfried Webers NICTA The Australian National University



Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class
Distributions

- Introduction to Statistical Machine Learning
- ©2013 Christfried Webers NICTA The Australian National



Maximum Margin Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

- Now for the bias b.
- Observe that for the support vectors,  $t_n y(\mathbf{x}_n) = 1$ .
- Therefore, we find

$$t_n\left(\sum_{m\in\mathcal{S}}a_m\,t_m\,k(\mathbf{x}_n,\mathbf{x}_m)+b\right)=1$$

and can use any support vector to calculate b.

• Numerically more stable: Multiply by  $t_n$ , observe  $t_n^2 = 1$ , and average over all support vectors.

$$b = \frac{1}{N_{\mathcal{S}}} \sum_{n \in \mathcal{S}} \left( t_n - \sum_{m \in \mathcal{S}} a_m t_m k(\mathbf{x}_n, \mathbf{x}_m) \right)$$

Maximum Margin Classifiers

Overlapping Class
Distributions

Relevance Vector Machines - Overview

 Express the error for maximum margin classifiers with the help of a function

$$E_{\infty}(z) = \begin{cases} 0, & z \ge 0\\ \infty, & z < 0 \end{cases}$$

enforcing the constraints as

$$\sum_{n=1}^{N} E_{\infty}(y(\mathbf{x}_n)t_n - 1)) + \lambda \|\mathbf{w}\|^2.$$



Classifiers

Overlapping Class Distributions

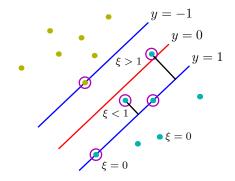
Relevance Vector Machines - Overview

- If the training data in the feature space are not linearly separable?
- Allow some data points to be on the 'wrong side' of the decision boundary.
- Increase a penalty with distance from the decision boundary.
- Assume, the penalty increases linearly with distance.
- Introduce slack variable  $\xi_n \geq 0$  for each data point n .

 $\xi_n \begin{cases} = 0, & \text{data point is correctly classified and} \\ & \text{on margin boundary or beyond} \\ < 1, & \text{data point is in correct margin} \\ = 1, & \text{data point is on the decision boundary} \\ > 1, & \text{data point is misclassified} \end{cases}$ 

• Introduce slack variable  $\xi_n \ge 0$  for each data point n.

$$\xi_n = \begin{cases} 0, & \text{data point is correctly classified and} \\ & \text{on margin boundary or beyond} \\ |t_n - y(\mathbf{x})|, & \text{otherwise} \end{cases}$$



Introduction to Statistical Machine Learning

> © 2013 Christfried Webers NICTA

The Australian National University



Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class Distributions





Overlapping Class Distributions

Constraints of the separable classification

$$t_n y(\mathbf{x}_n) \geq 1, \qquad n = 1, \dots, N$$

change now to

$$t_n y(\mathbf{x}_n) \geq 1 - \xi_n, \qquad n = 1, \dots, N$$

Minimise now

$$C\sum_{n=1}^{N}\xi_{n}+\frac{1}{2}\|\mathbf{w}\|^{2}$$

where C controls the trade-off between the slack variable penalty and the margin.

• Any misclassified point contributes  $\xi_n >$ , therefore  $\sum_{n=1}^N \xi_n$ is an upper bound on the number of misclassified points.

The Australian National University



Sparse Kernel Machines

Maximum Margin Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

• The Lagrangian with the Lagrange multipliers  $\mathbf{a} = (a_1, \dots, a_N)^T$  and  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_N)^T$  is now

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \mathbf{a}, \boldsymbol{\mu}) = \frac{1}{2} ||\mathbf{w}||^2 + C \sum_{n=1}^{N} \xi_n$$
$$- \sum_{n=1}^{N} a_n \{ t_n y(\mathbf{x}_n) - 1 + \xi_n \} - \sum_{n=1}^{N} \mu_n \xi_n.$$

• The KKT conditions for n = 1, ..., N are

$$a_n \ge 0$$

$$t_n y(\mathbf{x}_n) - 1 + \xi_n \ge 0$$

$$a_n(t_n y(\mathbf{x}_n) - 1 + \xi_n) = 0$$

$$\mu_n \ge 0$$

$$\xi_n \ge 0$$

$$\mu_n \xi_n = 0.$$

Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

ullet The dual Lagrangian after eliminating  ${f w}, b, {m \xi}$ , and  ${m \mu}$  is then

$$\widetilde{L}(\mathbf{a}) = \sum_{n=1}^{N} a_n - \frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\sum_{n=1}^{N} a_n t_n = 0$$

$$0 < a_n < C$$

$$n = 1, \dots, N.$$

 The only change from the separable case, is the box constraint via the parameter C.



Classifiers

Overlapping Class Distributions

Machines - Overview

Equivalent formulation by Schölkopf et. al. (2000)

$$\widetilde{L}(\mathbf{a}) = -\frac{1}{2} \sum_{n=1}^{N} \sum_{m=1}^{N} a_n a_m t_n t_m k(\mathbf{x}_n, \mathbf{x}_m)$$

$$\sum_{n=1}^{N} a_n t_n = 0$$

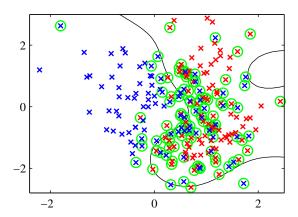
$$0 \le a_n \le 1/N$$

$$\sum_{n=1}^{N} a_n \ge \nu \qquad n = 1, \dots, N.$$

### where $\nu$ is both

- an upper bound on the fraction of margin errors, and,
- a lower bound on the fraction of support vectors.

# Overlapping Class distributions



The  $\nu$ -SVM algorithm using Gaussian kernels  $\exp(-\gamma \|\mathbf{x}-\mathbf{x}'\|^2)$  with  $\gamma=0.45$  applied to a nonseparable data set in two dimensions. Support vectors are indicated by circles.

Introduction to Statistical Machine Learning

© 2013 Christfried Webers NICTA The Australian National



Sparse Kernel Machines

Maximum Marg Classifiers

Overlapping Class Distributions

© 2013 Christfried Webers NICTA The Australian National University



- Sparse Kernel Machines
  - Maxımum Margın Classifiers
- Overlapping Class Distributions
- Relevance Vector Machines - Overview

- Output are decisions, not posterior probabilities.
- Extension to classification with more than two classes is problematic.
- There is a complexity parameter C (or  $\nu$ ) which must be found (e.g. via cross-validation).
- Predictions are expressed as linear combinations of kernel functions that are centered on the training points.
- Kernel matrix is required to be positive definite.



Maximum Margii Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview

Assume a Bayesian model as in Bayesian Linear
 Regression with the probability of the target t given by

$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}), \beta^{-1}),$$

and the linear model

$$y(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

with fixed nonlinear basis functions  $\phi(\mathbf{x})$  (including a constant term to accommodate for the bias).

• But now the prior over  ${\bf w}$  includes a different precision  $\alpha_i$  for each of the components of  ${\bf w}$ 

$$p(\mathbf{w} \mid \boldsymbol{\alpha}) = \prod_{i=1}^{M} \mathcal{N}(w_i \mid 0, \alpha_i^{-1}).$$



Relevance Vector Machines - Overview

When using this prior for the weights w

$$p(\mathbf{w} \mid \boldsymbol{\alpha}) = \prod_{i=1}^{M} \mathcal{N}(w_i \mid 0, \alpha_i^{-1})$$

and maximising the evidence with respect to the hyperparameters  $\alpha_i$ , a significant portion of the  $\alpha_i$  will go to infinity.

• The corresponding w<sub>i</sub> will therefore be concentrated at zero, and play no role for the prediction.



Maxımum Margın Classifiers

Overlapping Class
Distributions

- The Relevance Vector Machine iteratively calculates the  $\alpha_i$  and  $\beta$  from the training data by optimising a nonconvex function.
- Disadvantages
  - Nonconvex function means that the algorithm can get stuck in some local minimum.
  - Training times can be longer than for SVM.
- Advantages
  - As a Bayesian method, the parameters governing complexity and noise are determined from the input data. (Compare to SVM where we needed cross-validation).
  - Number of relevance vectors is much smaller than the number of support vectors in a SVM.
  - Therefore, prediction is much faster than with SVM.
  - Extends well to multiclass learning.

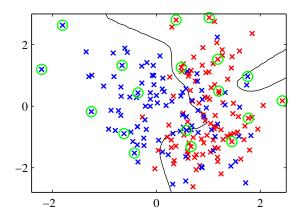




Classifiers

Overlapping Class Distributions

Relevance Vector Machines - Overview



Decision boundaries for a two-class problem found by a Relevance Vector Machine.





Maximum Margin

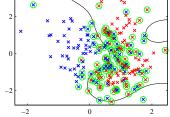
Overlapping Class

Relevance Vector Machines - Overview



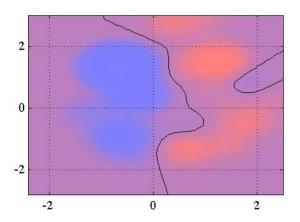
-2

Relevance Vector Machine



Support Vector Machine

### Relevance Vector Machines



Decision boundaries and posterior probability for a two-class problem found by a Relevance Vector Machine.

Introduction to Statistical Machine Learning

© 2013 Christfried Webers NICTA The Australian National University



Sparse Kernel Machines

Maximum Marg Classifiers

Overlapping Class
Distributions