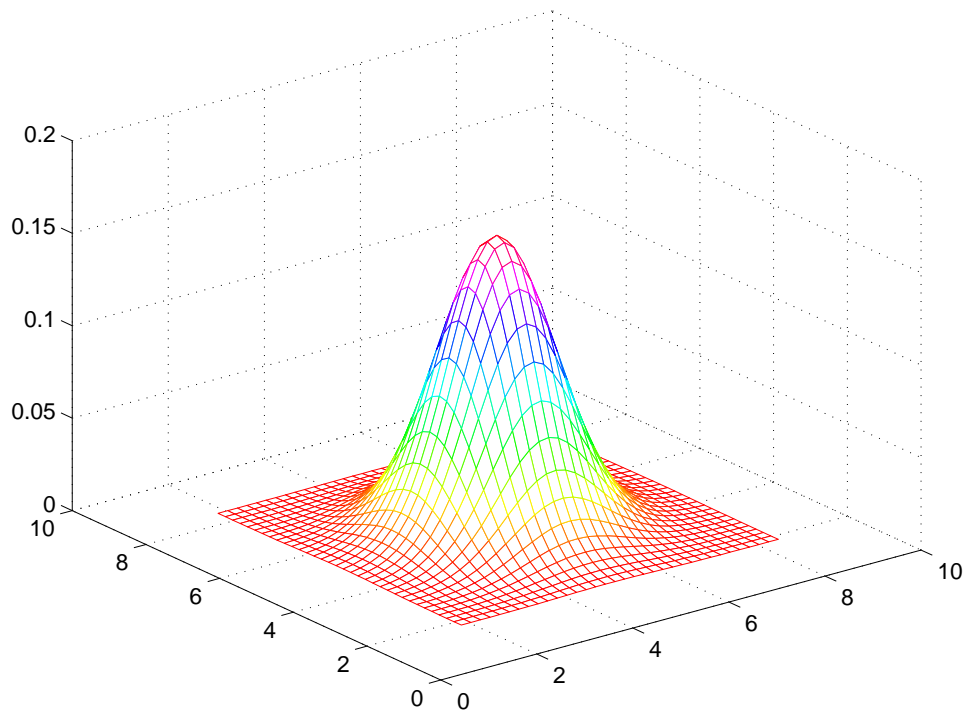


**University of Cambridge
Engineering Part IIB**

**Paper 4F10: Statistical Pattern
Processing**

Handout 9: Relevance Vector Machines



Bill Byrne
bill.byrne@eng.cam.ac.uk
Michaelmas 2012

Outline

The previous lectures have examined the use of Gaussian processes and basis function interpolation for prediction and classification. This was related to the support vector machines as both make use of kernel-like operations to yield non-linear classifiers.

These (the mean for Gaussian processes/basis interpolation) could be expressed in the form

$$\mu_y = \sum_{i=1}^n \alpha_i k(\mathbf{x}, \mathbf{x}_i)$$

One difference between these approaches and the SVM is that:

- **SVMs**: Lagrange multipliers α are **sparse**
- **GPs**: α **non-sparse**
- **basis interpolation**: in the general case **non-sparse**

Though SVMs are sparse there are some limitations:

- **SVMs are a discriminative function**: sometimes useful to have a discriminative model as this gives the class posteriors;
- **'C' must be set**: determines the error/margin play-off;

Would like to have some of the attributes of GPs/basis regression along with the sparseness of SVMs.

Basis Function Bayesian Regression

The underlying equations for basis function regression are

- **input-target relation**

$$y = \phi(\mathbf{x})' \mathbf{w} + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

- **weight distribution** is given by

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \sigma_w^2 \mathbf{I})$$

- **distribution of the training “targets”**

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \Phi \mathbf{w}, \sigma_n^2 \mathbf{I})$$

where

$$\Phi = \begin{bmatrix} \phi(\|\mathbf{x}_1 - \mathbf{x}_1\|) & \dots & \phi(\|\mathbf{x}_1 - \mathbf{x}_n\|) \\ \vdots & \ddots & \vdots \\ \phi(\|\mathbf{x}_n - \mathbf{x}_1\|) & \dots & \phi(\|\mathbf{x}_n - \mathbf{x}_n\|) \end{bmatrix} = \begin{bmatrix} \phi(\mathbf{x}_1)' \\ \vdots \\ \phi(\mathbf{x}_n)' \end{bmatrix}$$

Prediction is then based on

$$p(y | \mathbf{x}, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y; \phi(\mathbf{x})' \boldsymbol{\mu}_w, \phi(\mathbf{x})' \boldsymbol{\Sigma}_w \phi(\mathbf{x}) + \sigma_n^2)$$

where

$$p(\mathbf{w} | \mathbf{y}, \mathbf{X}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$$

and

$$\boldsymbol{\mu}_w = \frac{1}{\sigma_n^2} \left(\frac{1}{\sigma_n^2} \Phi' \Phi + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1} \Phi' \mathbf{y} = \frac{1}{\sigma_n^2} \boldsymbol{\Sigma}_w \Phi' \mathbf{y}$$

$$\boldsymbol{\Sigma}_w = \left(\frac{1}{\sigma_n^2} \Phi' \Phi + \frac{1}{\sigma_w^2} \mathbf{I} \right)^{-1}$$

All training observations contribute to the prediction - **not a sparse representation**

Gaussian Process Regression

The underlying equations for GPs are:

- **input-prediction relation:**

$$\begin{bmatrix} f(\mathbf{x}) \\ \mathbf{y} \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} k(\mathbf{x}, \mathbf{x}) & \mathbf{k}(\mathbf{x}, \mathbf{X})' \\ \mathbf{k}(\mathbf{x}, \mathbf{X}) & \mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I} \end{bmatrix}\right)$$

- **prediction-target relation**

$$y = f(\mathbf{x}) + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma_n^2)$$

The prediction mean can then be expressed as

$$p(y|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y; \mu_f, \sigma_f^2 + \sigma_n^2)$$

where

$$\begin{aligned} \mu_f &= \mathbf{k}(\mathbf{x}, \mathbf{X})' [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{y} \\ \sigma_f^2 &= k(\mathbf{x}, \mathbf{x}) - \mathbf{k}(\mathbf{x}, \mathbf{X})' [\mathbf{K}(\mathbf{X}, \mathbf{X}) + \sigma_n^2 \mathbf{I}]^{-1} \mathbf{k}(\mathbf{x}, \mathbf{X}) \end{aligned}$$

where $\mathbf{K}(\mathbf{X}, \mathbf{X})$ is determined by the **covariance function**

$$\mathbf{k}(\mathbf{x}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}, \mathbf{x}_1) \\ \vdots \\ k(\mathbf{x}, \mathbf{x}_n) \end{bmatrix}, \quad \mathbf{K}(\mathbf{X}, \mathbf{X}) = \begin{bmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & \dots & k(\mathbf{x}_n, \mathbf{x}_1) \\ \vdots & \ddots & \vdots \\ k(\mathbf{x}_1, \mathbf{x}_n) & \dots & k(\mathbf{x}_n, \mathbf{x}_n) \end{bmatrix}$$

Again all training observations contribute to the prediction - **not a sparse representation**

Relevance Vector Machine

The RVM is related to the basis regression form. However in order to create a more sparse representation the form of the weight prior is modified. Now

$$p(\mathbf{w}|\boldsymbol{\lambda}) = \prod_{i=1}^n \mathcal{N}(w_i; 0, \lambda_i^{-1}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \boldsymbol{\Lambda}^{-1})$$

Thus $\boldsymbol{\Lambda}$ is an inverse diagonal covariance matrix of the form

$$\boldsymbol{\Lambda} = \begin{bmatrix} \lambda_1 & \dots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \dots & \lambda_n \end{bmatrix}$$

Repeating the same process as before to get the weight posteriors yields the following values for the weight posteriors

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \mathcal{N}(\mathbf{w}; \boldsymbol{\mu}_w, \boldsymbol{\Sigma}_w)$$

where

$$\begin{aligned} \boldsymbol{\mu}_w &= \frac{1}{\sigma_n^2} \boldsymbol{\Sigma}_w \boldsymbol{\Phi}' \mathbf{y} \\ \boldsymbol{\Sigma}_w &= \left(\frac{1}{\sigma_n^2} \boldsymbol{\Phi}' \boldsymbol{\Phi} + \boldsymbol{\Lambda} \right)^{-1} \end{aligned}$$

Thus for the RVM there are the following parameters to train

- **noise variance:** σ_n^2 - the same as the Bayesian basis regression
- **weight-specific variance:** $\boldsymbol{\Lambda}$ - n parameters (matrix diagonal) rather than the single σ_w^2

RVM Prediction

Prediction with the RVM is very similar to Bayesian basis prediction

$$p(y|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \mathcal{N}(y; \phi(\mathbf{x})' \boldsymbol{\mu}_w, \phi(\mathbf{x})' \boldsymbol{\Sigma}_w \phi(\mathbf{x}) + \sigma_n^2)$$

Examining the just the mean of the prediction

$$\mu_y = \frac{1}{\sigma_n^2} \phi(\mathbf{x})' \left(\frac{1}{\sigma_n^2} \boldsymbol{\Phi}' \boldsymbol{\Phi} + \boldsymbol{\Lambda} \right)^{-1} \boldsymbol{\Phi}' \mathbf{y}$$

it is interesting to see the impact of varying λ_i on the prediction of the RVM.

- $\lambda_i \rightarrow \infty$ this means that element i of $\phi(\mathbf{x})$ will have no impact on the prediction.

Remembering that

$$\phi(\mathbf{x}) = \left[\phi(\|\mathbf{x} - \mathbf{x}_1\|) \quad \dots \quad \phi(\|\mathbf{x} - \mathbf{x}_n\|) \right]'$$

this means that training observation i will not alter the prediction - **it's not relevant to the prediction process** (see examples paper). This is related to the process of **automatic relevance detection**, and is the motivation for the name.

If a large number of the λ_i go to ∞ then there will be a **sparse prediction**.

In practice this is exactly what happens!

RVM Training

RVM training is based in the marginal likelihood

$$\begin{aligned} p(\mathbf{y}|\mathbf{X}) &= \int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}|\Lambda)d\mathbf{w} \\ &= \mathcal{N}(\mathbf{y}; \mathbf{0}, \Phi\Lambda\Phi' + \sigma_n^2\mathbf{I}) \end{aligned}$$

The values of Λ and σ_n^2 are then found by maximising this marginal likelihood.

The parameters may be estimated in the same fashion as the hyper-parameters of the Gaussian process using gradient descent, noting that

$$\Phi\Lambda\Phi' = \mathbf{K}(\mathbf{X}, \mathbf{X})$$

in the notation from the previous handout.

Alternatively **expectation maximisation** (EM) can be used. Here the value of the weight is treated as a continuous latent variable. This yields the **auxiliary function**

$$\begin{aligned} \mathcal{Q}(\boldsymbol{\theta}^{(k)}, \boldsymbol{\theta}^{(k+1)}) &= \int p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^{(k)}) \log \left(p(\mathbf{w}, \mathbf{y}|\mathbf{X}, \boldsymbol{\theta}^{(k+1)}) \right) d\mathbf{w} \\ &= \int p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^{(k)}) \log \left(p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \boldsymbol{\theta}^{(k+1)}) \right) d\mathbf{w} \\ &\quad + \int p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^{(k)}) \log \left(p(\mathbf{w}|\boldsymbol{\theta}^{(k+1)}) \right) d\mathbf{w} \end{aligned}$$

where

$$\boldsymbol{\theta} = \begin{bmatrix} \lambda_1 \\ \vdots \\ \lambda_n \\ \sigma_n^2 \end{bmatrix}$$

RVM Training (cont)

From the previous slides

$$p(\mathbf{y}|\mathbf{X}, \mathbf{w}, \boldsymbol{\theta}^{(k+1)}) = \mathcal{N}(\mathbf{y}; \Phi\mathbf{w}, \sigma_n^{(k+1)2})$$

$$p(\mathbf{w}|\boldsymbol{\theta}^{(k+1)}) = \mathcal{N}(\mathbf{w}; \mathbf{0}, \Lambda^{(k+1)-1})$$

$$p(\mathbf{w}|\mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^{(k)}) = \mathcal{N}(\mathbf{w}, \boldsymbol{\mu}_w^{(k)}, \boldsymbol{\Sigma}_w^{(k)})$$

Dealing with the two elements separately yields

- weights inverse variance

$$\lambda_i^{(k+1)} = \frac{1}{\sigma_{wi}^{(k)2} + \mu_{wi}^{(k)2}}$$

- noise variance

$$\sigma_n^{(k+1)2} = \frac{1}{n} \left((\mathbf{y} - \Phi\boldsymbol{\mu}_w^{(k)})'(\mathbf{y} - \Phi\boldsymbol{\mu}_w^{(k)}) + \text{tr}(\Phi'\Phi\boldsymbol{\Sigma}_w^{(k)}) \right)$$

This is an iterative process, at each iteration the marginal likelihood is **guaranteed** not to decrease.

For reference: it is possible to introduce priors on the hyper-parameters. Suitable priors are **Gamma** distributions

$$p(\Lambda) = \prod_{i=1}^n \text{Gamma}(\lambda_i|a, b)$$

$$p(\sigma_n^{-2}) = \text{Gamma}(\sigma_n^{-2}|c, d)$$

where

$$\text{Gamma}(\lambda|a, b) = \Gamma(a)^{-1} b^a \lambda^{a-1} \exp(-b\lambda)$$

$$\Gamma(a) = \int_0^\infty t^{a-1} \exp(-t) dt$$

RVM EM Training (details)

The estimation of the **weights inverse variance**

$$\Lambda^{(k+1)} = \arg \max_{\Lambda} \left\{ \int p(\mathbf{w} | \mathbf{X}, \mathbf{y}, \boldsymbol{\theta}^{(k)}) \log (p(\mathbf{w} | \Lambda)) d\mathbf{w} \right\}$$

Noting that the mean for $p(\mathbf{w} | \Lambda^{(k+1)})$ is constrained to be zero and Λ is diagonal yields the update formulae.

The estimation of the **noise variance** is given by

$$\sigma_n^{(k+1)2} = \arg \max_{\sigma} \left\{ \int \mathcal{N}(\mathbf{w}, \boldsymbol{\mu}_w^{(k)}, \boldsymbol{\Sigma}_w^{(k)}) \log (p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2)) d\mathbf{w} \right\}$$

The function of σ^2 to find may be written as

$$\log (p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2)) = -\frac{n}{2} \log (2\pi\sigma^2) - \frac{1}{2\sigma^2} (\mathbf{y} - \Phi\mathbf{w})' (\mathbf{y} - \Phi\mathbf{w})$$

Considering part of the second term on the RHS (expectation over \mathbf{w})

$$\mathcal{E} \{ (\mathbf{y} - \Phi\mathbf{w})' (\mathbf{y} - \Phi\mathbf{w}) \} = (\mathbf{y} - \Phi\boldsymbol{\mu}_w^{(k)})' (\mathbf{y} - \Phi\boldsymbol{\mu}_w^{(k)}) + \text{tr} (\Phi' \Phi \boldsymbol{\Sigma}_w^{(k)})$$

Noting this uses the equality

$$\begin{aligned} \mathcal{E} \{ \mathbf{w}' \Phi' \Phi \mathbf{w} \} &= \text{tr} (\Phi' \Phi \mathcal{E} \{ \mathbf{w} \mathbf{w}' \}) \\ &= \text{tr} \left(\Phi' \Phi \left(\boldsymbol{\Sigma}_w^{(k)} + \boldsymbol{\mu}_w^{(k)} \boldsymbol{\mu}_w^{(k)'} \right) \right) \\ &= \text{tr} \left(\Phi' \Phi \boldsymbol{\Sigma}_w^{(k)} \right) + \boldsymbol{\mu}_w^{(k)'} \Phi' \Phi \boldsymbol{\mu}_w^{(k)} \end{aligned}$$

Note: the update expression can be further simplified to yield

$$\sigma_n^{(k+1)2} = \frac{1}{n} \left((\mathbf{y} - \Phi\boldsymbol{\mu}_w^{(k)})' (\mathbf{y} - \Phi\boldsymbol{\mu}_w^{(k)}) + \sigma_n^{(k)2} \sum_{i=1}^n (1 - \lambda_i^{(k)} \sigma_{wi}^{(k)2}) \right)$$

RVM Classification

RVM classification uses exactly the same form as classification with Gaussian processes. Again the regression “prediction” distribution can be mapped using either:

- **logistic regression**: here

$$P(y = +1|\mathbf{x}, \mathbf{w}) = \frac{1}{1 + \exp(-\mathbf{w}'\phi(\mathbf{x}))}$$

- **probit regression**: again considering linear regression

$$P(y = +1|\mathbf{x}, \mathbf{w}) = \int_{-\infty}^{\mathbf{w}'\phi(\mathbf{x})} \mathcal{N}(z; 0, 1) dz$$

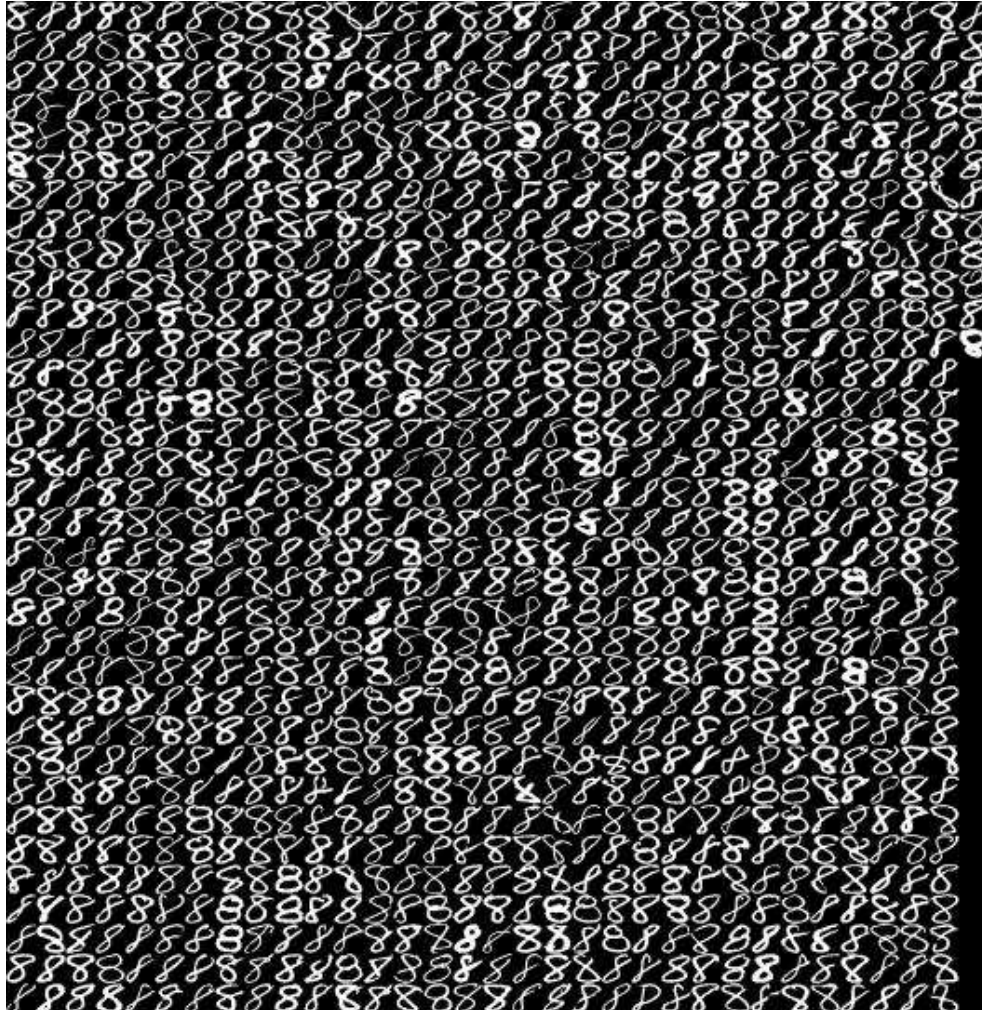
Classification is then based on

$$P(y = +1|\mathbf{x}, \mathbf{X}, \mathbf{y}) = \int P(y = +1|\mathbf{x}, \mathbf{w}) p(\mathbf{w}|\mathbf{X}, \mathbf{y}) d\mathbf{w}$$

Laplace’s approximation is again often used to handle the non-Gaussian nature that results from the mapping.

RVM Example Application - USPS digits

Recognition of handwritten digits - standard task - all the eights shown below



Train a set of **1-v-rest** classifiers using either an SVM or RVM with Gaussian kernels.

| Classifier | Training | | Number SVs/RVs | Errors % |
|------------|----------|-----|-------------------|-------------|
| | n | d | | |
| SVM | 7291 | 256 | 2540 | 4.4 |
| RVM | | | 316 | 5.1 |

RVM Summary

Various aspects of RVMs examined:

- relationship to basis function regression
- a sparse representation using the weight prior
- training using EM
- classification in same fashion as Gaussian processes