Machine Learning

Kernel Methods

[1]. Gaussian Processes (GP's)

[2]. Relevant Vector Machines (RVM's)

Lesson 7

Linear regression example

Assume M basis functions:

$$\phi(x) = (\phi_1(x), \phi_2(x), \dots, \phi_M(x))$$

- Linear model: $y(x) = w^T \phi(x) = \sum_{k=1}^m w_k \phi_k(x)$ where **w** is a M-dimensional weight vector
- An isotropic Gaussian prior over w

$$p(w \mid a) = N(w \mid 0, a^{-1}I)$$

- For any w obtain a particular function y(x).
- Intuitively, we take a probability distribution over function y(x)

$$y = (y(x_1), y(x_2), \dots, y(x_N))$$

• Linear model: $y = \Phi W$

where $\Phi = [\Phi_{ik} = \Phi_k(x_i)]$ the **design matrix**.

Remark: y is a linear combination of Gaussians

• **y** is also Gaussian p(y) = N(y|0,K)

$$E[y] = E[\Phi w] = \Phi E[w] = 0$$

$$COV[y] = E[(y - E[y])(y - E[y])^{T}] = E[yy^{T}] =$$

$$= E[(\Phi w)(\Phi w)^{T}] = \Phi E[ww^{T}]\Phi^{T} = \frac{1}{a}\Phi\Phi^{T} = K$$

where K is a kernel matrix (Gram matrix) :

$$K_{ij} = k(x_i, x_j) = \frac{1}{a} \phi(x_i)^T \phi(x_j)$$

This is a particular example of Gaussian Process (GP) model

- ✓ A GP provides a probability distribution over functions y(x), such that the set of values y(x) evaluated at (x₁, ..., x_N} jointly have a Gaussian distribution
- ✓ A **GP** provides a **stochastic process y(x)** that gives the joint distribution for any set of values $(y(x_1),...,y(x_N))$
- ✓ In 2-dimensional input case GP is known as Gaussian Random Field

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML7 (5)

- In Gaussian Process the joint of y's is specified completely by second-order statistics (mean and variance)
- Usually, we don't have any knowledge of mean and thus we set it to zero (0). (equivalent to zero-mean prior of w : p(w|a)=N(w|0,a⁻¹I))
- Thus, GP is then defined by giving the covariance of y(x), given by the kernel function

$$E[y(x_i)y(x_j)] = k(x_i, x_j)$$

Gaussian Processes for Regression

- Training set: $D = \{(x_1, t_1), ..., (x_N, t_N)\}$ $y(x) \sim GP$
- Generative model for targets:

$$t_i = y_i + \mathcal{E}_i$$

• where:

$$y_i = y(x_i)$$

$$\varepsilon_i \sim N(0, \beta^{-1})$$

()

random noise variable (β : precision of noise)

• Therefore: $p(t_i | y_i) = N(t_i | y_i, \beta^{-1})$

Join distribution of N target values T=(t₁, t₂,...,t_N)

$$p(T | Y) = N(T | Y, \beta^{-1}I_N) \qquad I_N : \text{NxN unit matrix}$$

where **Y=(y_1, y_2, ..., y_N)**

 Assuming that Y is a GPs, p(Y) is a zero-mean Gaussian with kernel (covariance) matrix K:

$$p(Y) = N(Y \mid 0, K)$$

Kernel function K_{ij} = k(x_i, x_j) is chosen to express the correlation level of the corresponding values y(x_i), y(x_j)

Marginal distribution of target values $T = (t_1, t_2, ..., t_N)$ $p(T) = \int p(T | Y) p(Y) dY$

2.3. The Gaussian Distribution 93

Using the properties of Gaussians we take that

Marginal and Conditional Gaussians

Given a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \tag{2.113}$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})$$
(2.114)

the marginal distribution of ${\bf y}$ and the conditional distribution of ${\bf x}$ given ${\bf y}$ are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
(2.115)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y}-\mathbf{b})+\mathbf{\Lambda}\boldsymbol{\mu}\}, \boldsymbol{\Sigma})$$
 (2.116)

where

$$\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}.$$
 (2.117)

Marginal distribution of target values $T = (t_1, t_2, ..., t_N)$ $p(T) = \int p(T | Y) p(Y) dY$ 2.3. The Gaussian Distribution 93 Using the properties of Marginal and Conditional Gaussians Given a marginal Gaussian distribution for x and a conditional Gaussian distri-Gaussians we take that bution for y given x in the form p(Y) = N(Y | 0, K) $p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\mu, \Lambda^{-1})$ (2.113) $p(\boldsymbol{T} | \boldsymbol{Y}) = N(\boldsymbol{Y} | \boldsymbol{Y}, \boldsymbol{\beta}^{-1} \boldsymbol{I}_{N}) p(\boldsymbol{y} | \boldsymbol{x}) = \mathcal{N}(\boldsymbol{y} | \boldsymbol{A}\boldsymbol{x} + \boldsymbol{b}, \boldsymbol{L}^{-1})$ (2.114)the marginal distribution of y and the conditional distribution of x given y are $p(T) = N(T \mid 0, C)$ $p(y) = \mathcal{N}(y \mid A\mu + b, L^{-1} + A\Lambda^{-1}A^{T})$ $p(x \mid y) = \mathcal{N}(x \mid \Sigma\{A^{T}L(y - b) + A\mu\}, \Sigma)$ where $p(x \mid y) = \mathcal{N}(x \mid \Sigma\{A^{T}L(y - b) + A\mu\}, \Sigma)$ (2.115)(2.116) $\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}.$ where (2.117) $[C]_{ii} = C(x_i, x_i) = k(x_i, x_i) + \beta^{-1}\delta_{ii}$

Regression Problem: Make prediction

 Predict the target t_{N+1} for a new input x_{N+1} given the target values T

$$p(t_{N+1} \mid T) \quad ?$$

• Joint distribution of $T_{N+1} = (t_1, ..., t_N, t_{N+1})$

$$p(T_{N+1}) = N(T_{N+1} | 0, C_{N+1})$$

 C_{N+1} covariance matrix is of size (N+1) x (N+1) with elements from kernel function on N+1 inputs

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (11)

• **Target:** find the conditional distribution $p(t_{N+1} | T)$

Solution:

Partition the covariance matrix

$$C_{N+1} = \begin{pmatrix} C & k \\ k^T & c \end{pmatrix}$$

where

- covariance matrix C is of size N x N
- k is a vector with elements $k(x_i, x_{N+1})$, i=1, ..., N

$$-c = k(x_{N+1}, x_{N+1}) + \beta^{-1}$$

Exploit the Gaussian properties

$$\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}. \tag{2.65}$$

We also define corresponding partitions of the mean vector μ given by

$$\boldsymbol{\mu} = \begin{pmatrix} \boldsymbol{\mu}_a \\ \boldsymbol{\mu}_b \end{pmatrix} \tag{2.66}$$

and of the covariance matrix Σ given by

$$\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}.$$
(2.67)

From these we obtain the following expressions for the mean and covariance of the conditional distribution $p(\mathbf{x}_a | \mathbf{x}_b)$

$$\mu_{a|b} = \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_b - \mu_b)$$
(2.81)

$$\Sigma_{a|b} = \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}.$$
(2.82)

Exploit the Gaussian properties

 $\mathbf{x} = \begin{pmatrix} \mathbf{x}_a \\ \mathbf{x}_b \end{pmatrix}$.

$$T_{N+1} = \begin{pmatrix} t_{N+1} \\ T \end{pmatrix}$$

We also define corresponding partitions of the mean vecto $\mu = \begin{pmatrix} \mu_a \\ \mu_b \end{pmatrix} \qquad \mu_{N+1} = \begin{pmatrix} 0 \\ \mu_N \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$ and of the covariance matrix Σ given by $\Sigma = \begin{pmatrix} \Sigma_{aa} & \Sigma_{ab} \\ \Sigma_{ba} & \Sigma_{bb} \end{pmatrix}. \qquad C_{N+1} = \begin{pmatrix} c & k^T \\ k & C \end{pmatrix}$

From these we obtain the following expressions for the mean and covariance of the conditional distribution $p(\mathbf{x}_a | \mathbf{x}_b)$

$$\begin{split} \mu_{a|b} &= \mu_a + \Sigma_{ab} \Sigma_{bb}^{-1} (\mathbf{x}_b - \mu_b) \\ \Sigma_{a|b} &= \Sigma_{aa} - \Sigma_{ab} \Sigma_{bb}^{-1} \Sigma_{ba}. \end{split}$$

$$m(x_{N+1}) = k^{T} C^{-1} T$$

$$\sigma^{2}(x_{N+1}) = c - k^{T} C^{-1} k$$

Gaussian Process regression model

$$p(t_{N+1} | T) = N(t_{N+1} | m(x_{N+1}), \sigma_{N+1}^2(x_{N+1}))$$

$$m(x_{N+1}) = k^T C^{-1} T$$

$$\sigma^2(x_{N+1}) = c - k^T C^{-1} k$$

$$c = k(x_{N+1}, x_{N+1}) + \beta^{-1}$$

$$k = (k(x_1, x_{N+1}), \dots, k(x_N, x_{N+1}))$$

[1]. Useful remark

• A necessary constraint is that the covariance matrix $\begin{bmatrix} C \end{bmatrix}_{ii} = C(x_i, x_i) = k(x_i, x_i) + \beta^{-1} \delta_{ii}$

must be **positive definite**.

- If λ_i is an eigenvalue of kernel matrix K, then matrix C will have as eigenvalue $\lambda_i + \beta^{-1}$.
- Thus, it is sufficient K be positive semidefinite so that λ ≥ 0.
- Need for constructing valid kernels...

296 6. KERNEL METHODS

Techniques for Constructing New Kernels.

Given valid kernels $k_1(\mathbf{x}, \mathbf{x}')$ and $k_2(\mathbf{x}, \mathbf{x}')$, the following new kernels will also be valid:

$$k(\mathbf{x}, \mathbf{x}') = ck_1(\mathbf{x}, \mathbf{x}') \tag{6.13}$$

$$k(\mathbf{x}, \mathbf{x}') = f(\mathbf{x})k_1(\mathbf{x}, \mathbf{x}')f(\mathbf{x}')$$
(6.14)

$$k(\mathbf{x}, \mathbf{x}') = q(k_1(\mathbf{x}, \mathbf{x}')) \tag{6.15}$$

$$k(\mathbf{x}, \mathbf{x}') = \exp\left(k_1(\mathbf{x}, \mathbf{x}')\right) \tag{6.16}$$

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}') + k_2(\mathbf{x}, \mathbf{x}')$$
 (6.17)

$$k(\mathbf{x}, \mathbf{x}') = k_1(\mathbf{x}, \mathbf{x}')k_2(\mathbf{x}, \mathbf{x}')$$
 (6.18)

$$k(\mathbf{x}, \mathbf{x}') = k_3(\phi(\mathbf{x}), \phi(\mathbf{x}'))$$
 (6.19)

$$\mathbf{x}(\mathbf{x},\mathbf{x}') = \mathbf{x}^{\mathrm{T}}\mathbf{A}\mathbf{x}' \tag{6.20}$$

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a) + k_b(\mathbf{x}_b, \mathbf{x}'_b)$$
(6.21)

$$k(\mathbf{x}, \mathbf{x}') = k_a(\mathbf{x}_a, \mathbf{x}'_a)k_b(\mathbf{x}_b, \mathbf{x}'_b)$$
(6.22)

where c > 0 is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(\mathbf{x})$ is a function from \mathbf{x} to \mathbb{R}^M , $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M , \mathbf{A} is a symmetric positive semidefinite matrix, \mathbf{x}_a and \mathbf{x}_b are variables (not necessarily disjoint) with $\mathbf{x} = (\mathbf{x}_a, \mathbf{x}_b)$, and k_a and k_b are valid kernel functions over their respective spaces.

Equipped with these properties, we can now embark on the construction of more complex kernels appropriate to specific applications. We require that the kernel $k(\mathbf{x}, \mathbf{x}')$ be symmetric and positive semidefinite and that it expresses the appropriate form of similarity between x and x' according to the intended application. Here we consider a few common examples of kernel functions. For a more extensive discussion of 'kernel engineering', see Shawe-Taylor and Cristianini (2004).

[2]. Useful remark

$$m(x_{N+1}) = k^T C^{-1} T = \sum_{i=1}^N a_i k(x_i, x_{N+1})$$

where
$$a_i = \begin{bmatrix} C^{-1}T \end{bmatrix}_i$$

- If kernel k(. , .) is Gaussian, then we obtain a radial basis function (RBF) network model.
- Note that inversion of matrix C requires O(N³) computational cost.

Learning the Gaussian process

Use a parametric (θ) kernel function in the covariance matrix. An example:

$$k(x_{i}, x_{j}) = \theta_{0} e^{-\frac{\theta_{1}}{2} \|x_{i} - x_{j}\|^{2}} + \theta_{2} + \theta_{3} x_{i}^{T} x_{j}$$

 $\Theta = \{ \theta_0, \theta_1, \theta_2, \theta_3 \}$ is the set of the unknown kernel parameters.

These can be **estimated** using training examples.



Figure 6.5 Samples from a Gaussian process prior defined by the covariance function (6.63). The title above each plot denotes $(\theta_{n}, \theta_{1}, \theta_{2}, \theta_{3})$.

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (19)

Learning the Gaussian process

- Learning: fit the Gaussian process to the data.
- Use the log likelihood function as a measure.

$$\ln p(T \mid \theta) = -\frac{N}{2} \ln(2\pi) - \frac{1}{2} \ln|C| - \frac{1}{2} T^T C^{-1} T$$

• Non-convex function => a lot of local maxima

$$\hat{\theta} = \arg \max_{\theta} \{ \ln p(T \mid \theta) \}$$

Calculating the derivatives

$$\ln p(T \mid \theta) = -\frac{1}{2} \ln |C| - \frac{1}{2} T^T C^{-1} T$$

• Useful relations:

Similarly

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{A}\mathbf{B}) = \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}}.$$
 (C.20)

The derivative of the inverse of a matrix can be expressed as

$$\frac{\partial}{\partial x} \left(\mathbf{A}^{-1} \right) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \tag{C.21}$$

as can be shown by differentiating the equation $A^{-1}A = I$ using (C.20) and then right multiplying by A^{-1} . Also

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \operatorname{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$
(C.22)

• Let $\theta = (\theta_i)$ be the set of parameters. Then

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – **ML7** (21)

Calculating the derivatives

$$\ln p(T \mid \theta) = -\frac{1}{2} \ln |C| - \frac{1}{2} T^{T} C^{-1} T$$

698

C. PROPERTIES OF MATRICES

Useful relations:

Similarly

$$\frac{\partial}{\partial \mathbf{x}} (\mathbf{AB}) = \frac{\partial \mathbf{A}}{\partial \mathbf{x}} \mathbf{B} + \mathbf{A} \frac{\partial \mathbf{B}}{\partial \mathbf{x}}.$$
 (C.20)

The derivative of the inverse of a matrix can be expressed as

$$\frac{\partial}{\partial x} \left(\mathbf{A}^{-1} \right) = -\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \mathbf{A}^{-1} \tag{C.21}$$

as can be shown by differentiating the equation $A^{-1}A =$ using (C.20) and then right multiplying by A^{-1} . Also

$$\frac{\partial}{\partial x} \ln |\mathbf{A}| = \operatorname{Tr} \left(\mathbf{A}^{-1} \frac{\partial \mathbf{A}}{\partial x} \right)$$
(C.22)

• Let $\theta = (\theta_i)$ be the set of parameters. Then

 $\frac{\partial \ln p(T \mid \theta)}{\partial \theta_{i}} = -\frac{1}{2} Tr \left(C^{-1} \frac{\partial C}{\partial \theta_{i}} \right) + \frac{1}{2} T^{T} C^{-1} \frac{\partial C}{\partial \theta_{i}} C^{-1} T$

Gaussian processes for Classification

 Bayesian estimation: Model the posterior probability

 Adapt Gaussian Process model to classification problems

 Transform the output of GPs using an nonlinear activation function

Problem formulation

 \succ Consider a binary classification problem $t \in \{0,1\}$

Define a GP over a function a(x)

Use a *logistic sigmoid function* so as to obtain $y \in [0, 1]$

$$y = \sigma(a) = \frac{1}{1 + e^{-a}}$$

Probability distribution over target is Bernoulli

$$p(t \mid a) = \sigma(a)^t (1 - \sigma(a))^{1-t}$$

Prediction with Gaussian processes

• Training set $D = \{(x_1, t_1), \dots, (x_N, t_N)\}$ (input, output pairs)

> Introduce a GP over the a's $a = \{a(x_1), \dots, a(x_N)\}$

$$p(a) = N(a \mid 0, C)$$

$$C = K + \beta^{-1}I$$

$$K_{ij} = k(x_i, x_j)$$
• Generative (or transformation) scheme



Prediction with Gaussian processes

- Test point x_{N+1} with unknown class t_{N+1}
- Goal is to determine the posterior

$$p(t_{N+1} \mid T) \qquad T = (t_1, \dots, t_N)$$

> Joint distribution of $a_{N+1} = (a_1, ..., a_N, a_{N+1})$ $a_{N+1} = \{a(x_1), ..., a(x_N), a(x_{N+1})\}$

$$p(a_{N+1}) = N(a_{N+1} | 0, C_{N+1})$$

Prediction with Gaussian processes

• For 2-class problems we want to predict the $p(t_{N+1} = 1 | T) \quad p(t_{N+1} = 0 | T_N) = 1 - p(t_{N+1} = 1 | T_N)$

$$p(t_{N+1} = 1 | T) =$$

= $\int p(t_{N+1} = 1 | a(x_{N+1})) p(a(x_{N+1}) | T) da(x_{N+1})$

where
$$p(t_{N+1} = 1 | a(x_{N+1})) = \sigma(a(x_{N+1}))$$

 $p(a(x_{N+1}) | T) = N(a(x_{N+1}) | m(x_{N+1}, \sigma_{N+1}^2(x_{N+1})))$
Obtaining from GP regression

Remember GP regression

$$p(t_{N+1} | T) = N(t_{N+1} | m(x_{N+1}), \sigma_{N+1}^2(x_{N+1}))$$

$$m(x_{N+1}) = k^T C^{-1} T$$

$$\sigma^2(x_{N+1}) = c - k^T C^{-1} k$$

$$c = k(x_{N+1}, x_{N+1}) + \beta^{-1}$$

$$k = (k(x_1, x_{N+1}), \dots, k(x_N, x_{N+1}))$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (28)

Approximation methods for integral calculation

$$p(t_{N+1} = 1 | T) = \int \sigma(a(x_{N+1})) N(a(x_{N+1}) | m_{N+1}, \sigma_{N+1}^2) da(x_{N+1}) = E_{N(m_{N+1}, \sigma_{N+1}^2)} [\sigma(a(x_{N+1}))]$$

- Laplace approximation [Barber & Williams]
- Variational methods [Gibbs & MacKay]
- Expectation-Propagation [Minka & Ghahramani]
- MCMC sampling [Neal]

Summary

- GP's provide a **structured method of model** and parameter selection.
- The key ingredient of a GP is the covariance function; a recipe to construct covariance matrices.
- Learning takes the form of setting the hyperparameters, using the **marginal likelihood**.
- GP's can be used for regression or classification. However require approximate inference techniques.

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (30)

Further reading on Gaussian Processes

Carl Edward Rasmussen and Chris Williams, MIT Press, 2006

Many more topics and code:

http://www.gaussianprocess.org/

Gaussian Processes for Machine Learning



Carl Edward Rasmussen and Christopher K. I. Williams

Relevance Vector Machine – RVM

Michael E. Tipping, "Sparse Bayesian Learning and the Relevance Vector Machine", JMLR, 2001

- The relevance vector machine (RVM) is a Bayesian sparse kernel method for regression and classification.
- It covers many applications
- Solves problems with the SVM

Support Vector Machines (SVM)

- A non-probabilistic decision machine: Returns point estimate for decision
- Makes decisions based on the function:

$$y(x,w) = \sum_{i=1}^{N} w_i K(x,x_i) + w_0$$

where K is the kernel function

 Attempts to minimize the error while simultaneously maximize the margin between the two classes

SVM "Problems"

> Non-probabilistic predictions.

- Requires estimation of error/margin trade-off parameters
- The kernel function K(x, x_i) must satisfy mercer's condition:
 - K must be a positive definite function

$$\iint f(x)f(y)K(x,y)dxdy > 0 \quad (\forall f \in L_2)$$
$$\int_{-\infty}^{+\infty} |f(x)|^2 dx < \infty$$

Relevance Vector Machines – RVM's

- Apply **Bayesian treatment** to SVM
- The kernel functions in RVM are treated simply as a set of basis functions without many restrictions imposed on SVM kernels
- Sparseness: Posterior distributions of the majority of weights are peaked around zero. Training vectors associated with the non-zero weights are the 'relevant vectors'.
- Uses significantly fewer kernel functions than SVM.

RVM for regression

- Training set $D = \{(x_1, t_1), \dots, (x_N, t_N)\}$
- Generative model

$$t = y(x) + \varepsilon = w^T \phi(x) + \varepsilon$$

- Assuming zero mean Gaussian noise $\mathcal{E} \sim N \Big(0, \beta^{-1} \Big)$
- We take the conditional distribution of target

$$p(t \mid x, w, \beta) = N(t \mid w^T \phi(x), \beta^{-1})$$

RVM for regression

• "Classical linear regression" assumes a linear combination with M nonlinear basis functions

$$y(x) = w^T \phi(x) = \sum_{i=1}^M w_i \phi_i(x)$$

• **RVM assumes** N kernels, one for each training example:

$$y(x) = \sum_{i=1}^{N} w_i k(x, x_i) + w_0$$

 It has the same structure of SVM, except that coefficients a_i are now denoted as w_i. • Likelihood function for the set of N target values

$$p(T \mid X, w, \beta) = \prod_{i=1}^{N} p(t_i \mid x_i, w, \beta) = N(T \mid \Phi w, \beta^{-1}I)$$

$$p(T \mid X, w, \beta) = \left(\frac{\beta}{2\pi}\right)^{\frac{N}{2}} \exp\left\{-\frac{\beta}{2} \left\|T - \Phi w\right\|^{2}\right\}$$

where
$$T = (t_1, t_2, ..., t_N)$$
 $w = (w_0, w_1, ..., w_N)$
 $\Phi = [\phi(x_1)\phi(x_2)\dots\phi(x_N)]$
 $\phi(x_i) = [k(x_1, x_i)k(x_2, x_i)\dots k(x_N, x_i)]$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (38)

ARD Prior: Gaussian prior for weights defining a separate hyperparameter a_i ∀ w_i:

$$p(w \mid a) = \prod_{i=0}^{N} N(w_i \mid 0, a_i^{-1}) = N(w \mid 0, A^{-1})$$
$$A = diag(a_1, \dots, a_N)$$

- a: is a vector of N+1 hyperparameters
- Introduce an hyperprior over a and precision $\boldsymbol{\beta}$

$$p(a) = \prod_{i=0}^{N} Gamma(a_i \mid \alpha, b)$$
$$p(\beta) = Gamma(\beta \mid c, d)$$

 Introduce sparse weights: integrate out the precision a and take the marginal weight prior

$$p(w_{i}) = \int p(w_{i} \mid a_{i})p(a_{i})da_{i} =$$

$$= \int N(w_{i} \mid 0, \alpha_{i}^{-1})Gamma(\alpha_{i} \mid a, b)d\alpha_{i} =$$

$$= \frac{b^{a}\Gamma\left(a + \frac{1}{2}\right)}{(2\pi)^{\frac{1}{2}}\Gamma(a)}\left(b + \frac{w_{i}^{2}}{2}\right)^{-\left(a + \frac{1}{2}\right)}$$
Student-t
distribution

Special case where *a=b=0 =>* Uninformative priors

$$p(w_i) \propto \frac{1}{|w_i|}$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (40)

Student-t distribution (Gosset 1908) $T \sim t_n$

•
$$E(T) = 0$$
 $VAR(T) = \frac{n}{n-2}$
• $\lim_{n \to \infty} f_n(t) = \frac{1}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$ $N(0,1)$

pdf:



$$p_{n}(t) = \frac{\Gamma\left(\frac{n+1}{2}\right)}{\sqrt{\pi}\sqrt{n}\Gamma\left(\frac{n}{2}\right)} \left(1 + \frac{t^{2}}{n}\right)^{-\frac{n+1}{2}} = \frac{1}{\sqrt{\pi}B\left(\frac{1}{2}, \frac{n}{2}\right)} \left(1 + \frac{t^{2}}{n}\right)^{-\frac{n+1}{2}}$$



Figure 6: LEFT: an example Gaussian prior $p(\mathbf{w}|\alpha)$ in two dimensions. RIGHT: the prior $p(\mathbf{w})$, where the hyperparameters have been integrated out to give a product of Student-*t* distributions. Note that the probability mass is concentrated both at the origin and along 'spines' where one of the two weights is zero.

Priors - graphical model





<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (44)





<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (46)

Inference – posterior over w

• The posterior over the weights w

$$p(w \mid T, a, \beta) = \frac{p(T \mid w, \beta)p(w \mid \alpha)}{\int p(T \mid w, \beta)p(w \mid \alpha)dw}$$

 Since all distributions are Gaussian, we can obtain analytical expression for the posterior pdf: product of two Gaussians

$$p(w | T, \alpha, \beta) \propto p(T | w, \beta) p(w | \alpha) =$$
$$= N(T | \Phi w, \beta^{-1}I) N(w | 0, A^{-1})$$

If we have a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
 (B.42)

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}+\mathbf{b},\mathbf{L}^{-1})$$
 (B.43)

then the marginal distribution of y, and the conditional distribution of x given y, are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
 (B.44)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y}-\mathbf{b})+\mathbf{\Lambda}\boldsymbol{\mu}\},\mathbf{\Sigma})$$
 (B.45)

where

$$\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}. \tag{B.46}$$

C. M. Bishop, "Pattern Recognition and Machine Learning", page 689

Machine Learning 2017 – Computer Science & Engineering, University of Ioannina – ML7 (48)

• Posterior distribution: product of two Gaussians

$$p(w | T, \alpha, \beta) \propto p(T | w, \beta) p(w | \alpha) =$$
$$= N(T | \Phi w, \beta^{-1}I) N(w | 0, A^{-1})$$

If we have a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \frac{N(w|0, A^{-1})}{N(y|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})} \frac{N(w|0, A^{-1})}{N(T | \Phi_{W}, \beta^{-1}I)^{B.43}}$$

then the marginal distribution of y, and the conditional distribution of x given y, are given by

$$p(w|T,X) = N(w|\mu,\Sigma)_{\Lambda^{-1}\Lambda^{\mathrm{T}}}$$
(B.44)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y}-\mathbf{b})+\mathbf{\Lambda}\boldsymbol{\mu}\},\mathbf{\Sigma})$$
 (B.45)

where

$$\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}.$$

C. M. Bishop, "Pattern Recognition and Machine Learning", page 689

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML7 (49)

 $\mu = \beta \Sigma \Phi^{T} T_{(A.46)}$ $\Sigma = \left(A + \beta \Phi^{T} \Phi\right)^{-1}$

• **Posterior distribution** of weights

$$p(w|T,a,\beta) = N(w|\mu,\Sigma)$$

$$\mu = \beta \Sigma \Phi^T T$$

$$\Sigma = \left(A + \beta \Phi^T \Phi\right)^{-1}$$

Estimation of model parameters

2 approaches

1st approach

• Maximum A-Posteriori (MAP) estimation problem

 $\ln p(w, \alpha, \beta | T) \propto \\ \propto \ln p(T | w, \beta) p(w | \alpha) p(\alpha | a, b) p(\beta | c, d)$

• Maximizing the log-likelihood (e.g. using EM algorithm)

$$\{\theta = \{w, \alpha, \beta\}\} = \max_{\theta} \ln p(w, \alpha, \beta | T)$$

Using **EM algorithm** for Maximizing the MAP log-likelihood

• Treat weights as hidden variables and maximize

$$E_{w|T,\alpha,\beta}\left[\ln p(T \mid w,\beta)p(w \mid \alpha)p(\alpha)p(\beta)\right]$$

• Update rules:

$$\alpha_{i} = \frac{1+2a}{E_{w|T,\alpha,\beta} [w_{i}^{2}] + 2b} = \frac{1+2a}{\Sigma_{ii} + \mu_{i}^{2} + 2b}$$
$$\left(\beta^{-1}\right)^{new} = \frac{\|T - \Phi\mu\|^{2} + (\beta^{-1})^{old} \sum_{i} \gamma_{i} + 2d}{N + 2c}$$

$$\gamma_i \equiv 1 - a_i \Sigma_i$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (52)

Estimation of model parameters

2nd approach

• Obtaining the marginal likelihood of targets

$$p(T \mid X, \alpha, \beta) = \int p(T \mid X, w, \beta) p(w \mid \alpha) dw$$

Maximum likelihood estimation problem

$$\{\theta = \{\alpha, \beta\}\} = \max_{\theta} \ln p(T \mid X, \alpha, \beta)$$

Marginal likelihood

$$p(T \mid \alpha, \beta) = \int p(T \mid w, \beta) p(w \mid \alpha) dw$$

If we have a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1})$$
 (B.42)

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}+\mathbf{b},\mathbf{L}^{-1})$$
 (B.43)

then the marginal distribution of y, and the conditional distribution of x given y, are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
 (B.44)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y}-\mathbf{b})+\mathbf{\Lambda}\boldsymbol{\mu}\},\mathbf{\Sigma})$$
 (B.45)

where

$$\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}. \tag{B.46}$$

C. M. Bishop, "Pattern Recognition and Machine Learning", page 689

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – **ML7** (54) 54

Marginal likelihood

$$p(T \mid \alpha, \beta) = \int p(T \mid w, \beta) p(w \mid \alpha) dw$$

If we have a marginal Gaussian distribution for x and a conditional Gaussian distribution for y given x in the form

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \frac{N(w|0, A^{-1})}{N(y|\mathbf{A}\mathbf{x} + \mathbf{b}, \mathbf{L}^{-1})} \frac{N(w|0, A^{-1})}{N(T | \Phi w, \beta^{-1}I)^{B.43}}$$

then the marginal distribution of y, and the conditional distribution of x given y, are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
 (B.44)

$${}^{p} p(T \mid \alpha, \beta) = N(T \mid 0, S)^{(B.45)}$$

where

 $S = \beta^{-1}I + \Phi A^{-1}\Phi^T$ (B.46)

55

C. M. Bishop, "Pattern Recognit

Maximize marginal log-likelihood

$$L = -\frac{1}{2}\ln|S| - \frac{1}{2}T^{T}S^{-1}T + \sum_{i=0}^{N}\ln Gamma(\alpha_{i} \mid a, b) + \ln p(\beta)$$

• Maximizing wrt ln(α) and log(β) and using $p(\ln a) = ap(a)$

$$L = -\frac{1}{2} \left\{ \ln |S| + T^T S^{-1} T \right\} + \sum_{i=0}^{N} \left(a \ln \alpha_i - b \alpha_i \right) + c \ln \beta - d\beta$$

where

$$S = \beta^{-1}I + \Phi A^{-1}\Phi^T$$

• Taking derivatives equal to zero:

$$\frac{\partial L}{\partial \ln \alpha_i} = \frac{1}{2} \left[1 - \alpha_i \left(\mu_i^2 + \Sigma_{ii} \right) \right] + a - b \alpha_i = 0 \Longrightarrow$$

$$\alpha_i^{new} = \frac{1+2a}{\mu_i^2 + \Sigma_{ii} + 2b}$$

$$\Sigma = \left(A + \beta \Phi^T \Phi\right)^{-1}$$

• Alternative rule using the quantities:

$$\gamma_i \equiv 1 - a_i \Sigma_{ii}$$

MacKay 1992

$$\alpha_i^{new} = \frac{\gamma_i + 2a}{\mu_i^2 + 2b}$$

faster convergence

• Alternative rule using the quantities:

$$\gamma_i \equiv 1 - a_i \Sigma_{ii}$$

$$\alpha_i^{new} = \frac{\gamma_i + 2a}{\mu_i^2 + 2b} \qquad \Sigma = (A + \beta \Phi^T \Phi)^{-1}$$

- Quantities *γ_i* ∈ [0,1] show how well-determined parameter *w* is by the data (*MacKay*, 1992).
 - If a_i is large (w_i doesn't fits the data) => $\Sigma_{ii} \approx \alpha_i^{-1}$ and thus $\gamma_i \approx 0$
 - If a_i is small (w_i fits the data) => $\gamma_i \approx 1$

• Taking derivatives equal to zero:

$$\frac{\partial L}{\partial \ln \beta} = \frac{1}{2} \left[\frac{N}{\beta} - \left\| T - \Phi \mu \right\|^2 - Tr \left(\Sigma \Phi^T \Phi \right) \right] + c - d\beta = 0 \Longrightarrow$$

• Using the fact that $Tr(\Sigma \Phi^T \Phi) = \beta^{-1} \sum_i \gamma_i$

$$\gamma_i \equiv 1 - a_i \Sigma_{ii}$$
$$\Sigma = \left(A + \beta \Phi^T \Phi\right)^{-1}$$

then:

$$\left(\beta^{-1}\right)^{new} = \frac{\left\|T - \Phi\mu\right\|^2 + 2d}{N - \sum_i \gamma_i + 2c}$$

Inference – Making predictions

 Given a new input x^{*} make a prediction of its target t^{*}. This is equivalent on estimating posterior distribution

$$p(t^* | T = \{t_1, \dots, t_N\}) = p(t^* | x^*, T, \hat{\alpha}, \hat{\beta})$$

• Marginalizing we take:

$$p(t^* | x^*, T, \hat{\alpha}, \hat{\beta}) = \int p(t^* | x^*, w, \hat{\beta}) p(w | T, \hat{\alpha}, \hat{\beta}) dw$$

Integration of two Gaussians

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (60)

$$p(t^* | x^*, T, a, \beta) = \int p(t^* | x^*, w, \hat{\beta}) p(w | T, \hat{a}, \hat{\beta}) dw$$

where

$$p(t^* | x^*, w, \beta) = N(t^* | w^T \phi(x^*), \beta^{-1})$$

$$p(w | T, \hat{\alpha}, \hat{\beta}) = N(w | \mu, \Sigma)$$
$$\mu = \hat{\beta} \Sigma \Phi^T T \qquad \Sigma = (\hat{A} + \hat{\beta} \Phi^T \Phi)^{-1}$$

<u>Machine Learning 2017</u> – Computer Science & Engineering, University of Ioannina – ML7 (61)

Predictive distribution (cont.)

$$p(t^* | x^*, T, a, \beta) = \int p(t^* | x^*, w, \hat{\beta}) p(w | T, \hat{a}, \hat{\beta}) dw$$

$$p(\mathbf{x}) = \mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Lambda}^{-1}) \tag{B.42}$$

$$p(\mathbf{y}|\mathbf{x}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\mathbf{x}+\mathbf{b},\mathbf{L}^{-1})$$
 (B.43)

then the marginal distribution of y, and the conditional distribution of x given y, are given by

$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}})$$
 (B.44)

$$p(\mathbf{x}|\mathbf{y}) = \mathcal{N}(\mathbf{x}|\mathbf{\Sigma}\{\mathbf{A}^{\mathrm{T}}\mathbf{L}(\mathbf{y}-\mathbf{b})+\mathbf{\Lambda}\boldsymbol{\mu}\},\mathbf{\Sigma})$$
 (B.45)

where

$$\Sigma = (\Lambda + \mathbf{A}^{\mathrm{T}} \mathbf{L} \mathbf{A})^{-1}.$$
 (B.46)

where:

Predictive distribution (cont.)

$$p(t^* \mid x^*, T, a, \beta) = \int p(t^* \mid x^*, w, \hat{\beta}) p(w \mid T, \hat{a}, \hat{\beta}) dw$$

$$p(x) = \mathcal{N}(x \mid \mu, \Lambda^{-1}) \qquad p(w \mid T, \hat{a}, \hat{\beta}) = N(w \mid \mu, \Sigma)$$

$$p(y \mid x) = \mathcal{N}(y \mid Ax + b, L^{-1}) \qquad (B.43)$$

$$p(t^* \mid x^*, w) = N(t^* \mid w^T \phi(x^*), \beta^{-1})$$

then the marginal distribution of y, and the conditio $p(t | x, w) = N(t | w \phi(x), \beta^{-1})$ given by

where
$$p(\mathbf{y}) = \mathcal{N}(\mathbf{y}|\mathbf{A}\boldsymbol{\mu} + \mathbf{b}, \mathbf{L}^{-1} + \mathbf{A}\boldsymbol{\Lambda}^{-1}\mathbf{A}^{\mathrm{T}}) \qquad (B.44)$$
$$p\left(t^{*} \mid x^{*}, T\right) = \mathcal{N}\left(t^{*} \mid \boldsymbol{\mu}^{T}\boldsymbol{\phi}(x^{*}), \boldsymbol{\sigma}^{2}(x^{*})\right)$$
$$\Sigma = (\mathbf{A} + \mathbf{A}^{*}\mathbf{L}\mathbf{A})^{-1} \qquad (B.46)$$
$$\sigma^{2}\left(x^{*}\right) = \hat{\boldsymbol{\beta}}^{-1} + \boldsymbol{\phi}^{T}\left(x^{*}\right)\boldsymbol{\Sigma}\boldsymbol{\phi}(x^{*})$$

where:

• Predictive distribution:

$$p(t^* | x^*, T) = N(t^* | \mu^T \phi(x^*), \sigma^2(x^*))$$

• Prediction:

$$y(x^*;\mu) = \mu^T \phi(x^*) = \sum_{i=1}^N \mu_i k(x_i, x^*)$$

• Variance:

$$\sigma^2(x^*) = \hat{\beta}^{-1} + \phi^T(x^*) \Sigma \phi(x^*)$$

noise variance + uncertainty on the prediction of the weights

RVM for classification

• Consider a binary classification problem $t \in \{0,1\}$

• Output is calculated using a logistic sigmoid function $y(x) = \sigma(w^T \phi(x))$

$$p(w \mid a) = \prod_{i=1}^{N} N(w_i \mid 0, a_i^{-1}) = N(w \mid 0, A^{-1})$$

• Obtain the posterior distribution of weights

Machine Learning 2017 - Computer Science & Engineering, University of Ioannina - ML7 (65)

RVM for classification

- Posterior distribution over w: $p(w|a) = N(w|0, A^{-1})$ $p(w|T, a) = \frac{p(T|w)p(w|a)}{p(T|a)}$
- Class conditional distribution of inputs

$$p(T | w) = \prod_{i=1}^{N} y_i^{t_i} (1 - y_i)^{1 - t_i} \qquad y_i = \sigma(w^T \phi(x_i))$$

Log-likelihood

$$\ln p(w | T, a) \propto \sum_{i=1}^{N} \{t_i \ln y_i + (1 - t_i) \ln(1 - y_i)\} - \frac{1}{2} w^T A w$$

• Maximization problem
$$y_i = \sigma(w^T \phi(x_i))$$

$$\max_{w} \left\{ \sum_{i=1}^{N} \left\{ t_i \ln y_i + (1 - t_i) \ln(1 - y_i) \right\} - \frac{1}{2} w^T A w \right\}$$

Use Newton-Raphson optimization scheme

$$w^{(new)} = w^{(old)} - H^{-1} \nabla E(w) \qquad H = \nabla \nabla E(w)$$

• We obtain:

$$\nabla E(w) = \Phi^T (T - Y) - Aw$$
$$\nabla \nabla E(w) = -(\Phi^T B \Phi + A) \qquad B = [b_i = y_i(1 - y_i)]$$