# RELEVANCE VECTOR MACHINES RVMs

## Matrix Identities

(C.7) $\quad (A+BD^{-1}C)^{-1} = A^{-1} - A^{-1}B(D+CA^{-1}B)^{-1}CA^{-1}$

(C.14) $\quad |I_N + AB^T| = |I_M + A^T B|$

(C.15) $\quad |I_N + ab^T| = 1 + a^T b$

## Gaussian Identities

$\left[\begin{array}{l}\text{(2.113)} \quad p(x) = \mathcal{N}(x \mid \mu, \Lambda^{-1}) \\ \text{(2.114)} \quad p(y|x) = \mathcal{N}(y \mid Ax+b, L^{-1}) \end{array}\right.$

$\longrightarrow$ (2.115) $\quad p(y) = \mathcal{N}(y \mid A\mu+b, L^{-1} + A\Lambda^{-1}A^T)$

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

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

## Limitations of SVMs

① posterior probabilities ?

② two-class $\Rightarrow$ multi-class

③ $C, \nu, \epsilon$ parameter selection by cross-validation

④ positive definite kernels

## RVM for regression

$t$ : real-valued target value

$x$ : input vector

Gaussian Noise : $p(t \mid x, w, \beta) = \mathcal{N}(t \mid y(x), \beta^{-1})$

$\qquad\qquad\qquad\qquad \underset{\text{mean}}{\Downarrow} \quad \underset{\text{noise precision}}{\Downarrow}$

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

for RVMs, we assume the following model

$$y(x) = \sum_{n=1}^{N} w_n k(x, x_n) + b$$

$$\Downarrow$$

as $\phi(x)$ in a linear model
no restriction to p.d kernels

use the matrix notation and assume i.i.d.

$X = \begin{bmatrix} \vdots \\ x_n^T \\ \vdots \end{bmatrix} \}N$

$t = \begin{bmatrix} \vdots \\ t_n \\ \vdots \end{bmatrix} \}N$

likelihood :

$$p(t \mid X, w, \beta) = \prod_{n=1}^{N} p(t_n \mid x_n, w, \beta)$$

Now introduce the prior on w

Key : each weight parameter $w_i$ has a separate hyperparameter $\alpha_i$

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

if $\alpha_i \to \infty \Rightarrow$ high precision, zero variance

$\qquad\qquad \Rightarrow w_i$ centered at the mean $_{=0}$

$\qquad\qquad \Rightarrow$ sparse model

from the likelihood and the prior we can write the posterior as

$$p(w \mid t, \underline{X}, \alpha, \beta) = \mathcal{N}(w \mid m, \Sigma)$$

where $\quad m = \Sigma \{ \Phi^T \beta t \}$
$$\Sigma = (A + \beta \Phi^T \Phi)^{-1}$$

$\Phi_{ni} = \phi_i(x_n)$
$A = \text{diag}(\alpha_i)$

(2.116)

this is obtained by applying (2.113) + (2.114) → (2.116) to

as (2.113) ⇒ $\quad p(w \mid \alpha) = \prod_{i=1}^{M} \mathcal{N}(w_i \mid 0, \alpha_i)$

as (2.114) ⇒ $\quad p(t \mid \underline{X}, w, \beta) = \prod_{n=1}^{N} \mathcal{N}(\underbrace{w^T \phi(x_n)}_{y(x_n)}, \beta^{-1})$

$\left( \text{Note that } \Phi = K \text{ for RVM} \right)$

The optimal weights $w^*$ is given by the mean of the posterior
$$w^* = m = \beta \Sigma \Phi^T t$$

The next step is to decide the hyperparameters $\alpha, \beta$.

We use "evidence approximation".

EA ① compute the marginal likelihood
$$p(t \mid \underline{X}, \alpha, \beta) = \int p(t \mid \underline{X}, w, \beta) \, p(w \mid \alpha) \, dw$$
$$= \int p(t, w \mid \underline{X}, \alpha, \beta) \, dw$$

---

$p(t \mid \underline{X}, w, \beta)$ and $p(w, \alpha)$ are Gaussians

Again, we don't need to do the integration explicitly for marginalization.
Just by observing and applying (2.113) + (2.114) → (2.115) we get

$$p(t \mid \underline{X}, \alpha, \beta) = \int p(t \mid \underline{X}, w, \beta) \, p(w \mid \alpha) \, dw$$
$$= \mathcal{N}(t \mid 0, C)$$
$$C = \beta^{-1} I + \Phi A^{-1} \Phi^T \qquad A = \text{diag}(\alpha_i)$$

EA ② maximize the logarithm of the marginal likelihood w.r.t. $\alpha$ and $\beta$

$$\ln p(t \mid \underline{X}, \alpha, \beta) = \ln \mathcal{N}(t \mid 0, C)$$
$$= -\frac{1}{2} \{ N \ln(2\pi) + \ln |C| + t^T C^{-1} t \}$$

It takes a couple of pages to write down the derivations. For now we write the results only, and go into the details later.

take the derivatives and set them to zero, we get

$$\alpha_i^{new} = \frac{\gamma_i}{m_i^2}$$

$$(\beta^{new})^{-1} = \frac{\| t - \Phi m \|^2}{N - \Sigma_i \gamma_i} \qquad \gamma_i = 1 - \alpha_i \Sigma_{ii}$$

$$\ln p(t \mid \overline{X}, \alpha, \beta) = -\frac{1}{2}\left\{ N \ln(2\pi) + \ln|C| + t^T C^{-1} t \right\}$$

$$\frac{\partial}{\partial \alpha_i}\left\{ \ln p(t \mid \overline{X}, \alpha, \beta) \right\} = \frac{\partial \ln|C|}{\partial \alpha_i} + t^T \frac{\partial C^{-1}}{\partial \alpha_i} t$$

① 
$$\ln|C| = \ln\left|\beta^{-1} I + \Phi A^{-1} \Phi^T\right|$$
$$= \ln \beta^{-N}\left| I + \beta \Phi A^{-1} \Phi^T \right|$$
$$= -N \ln\beta + \ln\left| I + \beta \Phi A^{-1} \Phi^T \right|$$
$$= -N \ln\beta + \ln\left| I + \beta A^{-1} \Phi^T \Phi \right| \qquad (C.14)$$
$$= -N \ln\beta - \ln|A| + \ln\left| A + \beta \Phi^T \Phi \right|$$
$$= -N \ln\beta - \ln|A| + \ln\left| \Sigma^{-1} \right|$$

recall w posterior
$$\Sigma = (A + \beta \Phi^T \Phi)^{-1}$$

$$\frac{\partial \ln|C|}{\partial \alpha_i} = -\mathrm{Tr}\left( A^{-1} \frac{\partial A}{\partial \alpha_i} \right) + \mathrm{Tr}\left( \Sigma \frac{\partial \Sigma^{-1}}{\partial \alpha_i} \right) \qquad (C.22)$$

$$= -\frac{1}{\alpha_i} + \Sigma_{ii}$$

$$\frac{\partial A}{\partial \alpha_i} = \begin{bmatrix} 0 & & 0 \\ & {}^{\circ}1_{\circ} & \\ 0 & & 0 \end{bmatrix}$$

$$\frac{\partial \Sigma^{-1}}{\partial \alpha_i} \text{ is similar}$$

② 
$$C^{-1} = \left( \beta^{-1} I + \Phi A^{-1} \Phi^T \right)^{-1}$$
$$= \left( \beta^{-1}\left( I + \beta \Phi A^{-1} \Phi^T \right) \right)^{-1}$$
$$(C.7) \qquad = \beta\left( I + \beta \Phi A^{-1} \Phi^T \right)^{-1} = \beta\left\{ I - \Phi\left( \beta^{-1} A + \Phi^T \Phi \right)^{-1} \Phi^T \right\}$$
$$= \beta\left\{ I - \beta \Phi\left( A + \beta \Phi^T \Phi \right)^{-1} \Phi^T \right\} = \beta\left\{ I - \beta \Phi \Sigma \Phi^T \right\}$$

$$\frac{\partial t^T C^{-1} t}{\partial \alpha_i} = t^T \frac{\partial C^{-1}}{\partial \alpha_i} t$$

$$\frac{\partial C^{-1}}{\partial \alpha_i} = \frac{\partial}{\partial \alpha_i}\left\{ \beta\left( I - \beta \Phi \Sigma \Phi^T \right) \right\}$$
$$= -\beta^2 \Phi \frac{\partial \Sigma}{\partial \alpha_i} \Phi^T$$
$$= -\beta^2 \Phi \left( -\Sigma \begin{bmatrix} 0 & & \\ & {}^{\circ}1_{\circ} & \\ & & 0 \end{bmatrix} \Sigma \right) \Phi^T$$

by (C.21)
$$\frac{\partial}{\partial \alpha_i}(\Sigma)$$
$$= -\Sigma \frac{\partial \Sigma^{-1}}{\partial \alpha_i} \Sigma$$
$$\Downarrow$$
$$\begin{bmatrix} 0 & & \\ & {}^{\circ}1_{\circ} & \\ & & 0 \end{bmatrix}$$

$$\Rightarrow \frac{\partial t^T C^{-1} t}{\partial \alpha_i}$$
$$= t^T \beta \Phi \Sigma \begin{bmatrix} 0 & & \\ & {}^{\circ}1_{\circ} & \\ & & 0 \end{bmatrix} \Sigma \Phi^T \beta t$$
$$= m^T \begin{bmatrix} 0 & & \\ & {}^{\circ}1_{\circ} & \\ & & 0 \end{bmatrix} m = m_i^2$$

recall w posterior
$$m = \Sigma \Phi^T \beta t$$

combine the results of ① and ②
$$\frac{\partial \ln p(t \mid \overline{X}, \alpha, \beta)}{\partial \alpha_i} = \frac{1}{2\alpha_i} - \frac{1}{2}\Sigma_{ii} - \frac{1}{2}m_i^2$$

Set the derivative to zero
$$\alpha_i = \frac{1 - \alpha_i \Sigma_{ii}}{m_i^2}$$

$\alpha^*$ The update rule 
$$\alpha_i^{new} = \frac{\gamma_i}{m_i^2} \qquad \gamma_i = 1 - \alpha_i^{old} \Sigma_{ii}$$

We also need to compute the derivatives w.r.t. $\beta$

Again $\quad \ln p(t \mid \overline{X}, \alpha, \beta) = -\frac{1}{2}\{N \ln(2\pi) + \ln|C| + t^T C^{-1} t\}$

$$\frac{\partial}{\partial \beta} \ln p(t \mid \overline{X}, \alpha, \beta) = -\frac{1}{2}\left\{\frac{\partial}{\partial \beta}\ln|C| + \frac{\partial}{\partial \beta} t^T C^{-1} t\right\}$$

① recall $\quad \ln|C| = -N\ln\beta - \ln|A| + \ln|\Sigma^{-1}|$

$$\frac{\partial \ln|C|}{\partial \beta} = \frac{-N}{\beta} + \frac{\partial}{\partial \beta}\ln|\Sigma^{-1}| \qquad \Big| \begin{array}{l} A \text{ is irrelevant} \\ \text{to } \beta \end{array}$$

Write $\quad \dfrac{\partial}{\partial \beta}\ln|\Sigma^{-1}| = Tr\left(\Sigma \dfrac{\partial \Sigma^{-1}}{\partial \beta}\right) \qquad \Big| \begin{array}{l} \text{recall in posterior} \\ \Sigma^{-1} = (A+\beta\Phi^T\Phi) \end{array}$

(c.22)
$$= Tr(\Sigma \Phi^T \Phi)$$
$$= Tr(\Sigma \Phi^T \Phi + \beta^{-1}\Sigma A - \beta^{-1}\Sigma A)$$
$$= Tr\{\Sigma(\Phi^T\Phi\beta + A)\beta^{-1} - \beta^{-1}\Sigma A\}$$
$$= Tr\{(I - A\Sigma)\beta^{-1}\}$$

$$\frac{\partial \ln|C|}{\partial \beta} = -\frac{N}{\beta} + \frac{1}{\beta} Tr\{I - A\Sigma\} \qquad \Big| \begin{array}{l} \text{let} \\ \gamma_i = 1 - \alpha_i \Sigma_{ii} \end{array}$$

$$= -\frac{N}{\beta} + \frac{1}{\beta}\sum_{i=1}^{N}\gamma_i$$

② $\quad \dfrac{\partial t^T C^{-1} t}{\partial \beta} \qquad \Big| \text{using previous results}$

$$= \frac{\partial}{\partial \beta}\{t^T \beta(I - \beta\overline{\Phi}\Sigma\overline{\Phi}^T) t\}$$

$$= \frac{\partial}{\partial \beta}\{\beta t^T t - \beta^2 t\overline{\Phi}\Sigma\overline{\Phi}^T t\}$$

$$= t^T t - 2\beta t\overline{\Phi}\Sigma\overline{\Phi}^T t + \beta^2 t\overline{\Phi}\Sigma\overline{\Phi}^T\overline{\Phi}\Sigma\overline{\Phi}^T t$$

$$= \|t - \Phi m\|^2$$

$$\Big| m = \beta\Sigma\overline{\Phi}^T t \qquad \Big| \begin{array}{l} \frac{\partial \Sigma}{\partial \beta} \\ = -\Sigma\frac{\partial \Sigma^{-1}}{\partial \beta}\Sigma \\ = -\Sigma\overline{\Phi}^T\Phi\Sigma \end{array}$$

$$\frac{\partial}{\partial \beta}\ln p(t \mid \overline{X}, \alpha, \beta)$$
$$= \frac{1}{2}\frac{N}{\beta} - \frac{1}{2}\frac{1}{\beta}\sum_{i=1}^{N}\gamma_i - \frac{1}{2}\|t - \Phi m\|^2 = 0$$

The update rule

$\beta^*$
$$\frac{1}{\beta^{new}} = \frac{\|t - \Phi m\|^2}{N - \sum_{i=1}^{N}\gamma_i} \qquad \gamma_i = 1 - \alpha_i\Sigma_{ii}$$

Predictive Distribution

$$p(t \mid x, \underset{\sim}{X}, t\!t, \alpha^*, \beta^*)$$

$$= \int \underbrace{p(t \mid x, w, \beta^*)}_{\mathcal{N}(t \mid y(x), \beta^{-1})} \underbrace{p(w \mid \underset{\sim}{X}, t\!t, \alpha^*, \beta^*)}_{\mathcal{N}(w \mid m, \Sigma)} dw$$

$$= \mathcal{N}(t \mid m^T \phi(x), \sigma^2(x))$$
$$\uparrow$$
$$(\alpha^*, \beta^*)$$

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

$$m = \beta^* \Sigma \, \Phi^T t\!t$$