

# Relevance vector machines (RVMs) } Lecture 16

SVM drawbacks: not probabilistic, hard to generalize to  $K > 2$ , nuisance prms  $C, \epsilon$ .

[ RVMs are a Bayesian technique for regression & classification which yields sparse models. ]

## Regression RVMs

Recall that in regression,

$$p(t|\vec{x}, \vec{w}, \beta) = \mathcal{N}(t | \underbrace{y(\vec{x})}_{\vec{w}^T \cdot \vec{g}(\vec{x})}, \beta^{-1})$$

"  $\beta^{-1} = \sigma^2$   
[bias included]

In RVMs, we use

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

$M = N+1$  prms in total

similar to  $y(\vec{x}) = \sum_n (a_n - \hat{a}_n) k(\vec{x}, \vec{x}_n) + b$   
in SVMs

Training data:

$$X = \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \\ \vdots \\ \vec{x}_N \end{pmatrix}$$

$$\vec{t} = \underline{t_1, t_2, \dots, t_N}$$

$$\text{Then } \mathcal{J} = p(\vec{t}|X, \vec{w}, b) = \prod_{n=1}^N p(t_n | \vec{x}_n, \vec{w}, \beta^{-1})$$

Introduce  $p(\vec{w} | \vec{d}) = \prod_{i=1}^M \mathcal{N}(w_i | 0, d_i^{-1})$

prior

separate  $d_i$  for each  $w_i$

$$\vec{d} = \overbrace{d_1 \dots d_M}$$

Then  $p(\vec{w} | \vec{t}, X, \vec{d}, \beta) = \mathcal{N}(\vec{w} | \vec{m}, \Sigma)$  ← posterior for  $\vec{w}$

$$\begin{cases} \vec{m} = \beta \Sigma \Phi^T \vec{t}, \\ \Sigma^{-1} = A + \beta \Phi^T \Phi \end{cases} \Leftrightarrow (3.53), (3.54)$$

$$A = \begin{pmatrix} d_1 & & 0 \\ & \ddots & \\ 0 & & d_M \end{pmatrix}$$

For RVMS,  $\Phi \Rightarrow K$ , with elements  $k(\vec{x}_n, \vec{x}_m)$   
 $(N+1) \times (N+1)$  kernel matrix  
 (including the bias)

$$y(\vec{x}) = \sum_{n=1}^N w_n k(\vec{x}, \vec{x}_n) + \underbrace{b}_{w_{N+1}} = \sum_{n=1}^{N+1} w_n k(\vec{x}, \vec{x}_n)$$

" 1,  $\forall \vec{x}$

Evidence approximation:

maximize  $p(\vec{t} | X, \vec{d}, \beta) = \underbrace{\int p(\vec{t} | X, \vec{w}, \beta)}_{\text{likelihood, gaussian}} \underbrace{p(\vec{w} | \vec{d})}_{\text{prior, gaussian}} d\vec{w}$

$$p(\vec{t} | X, \vec{d}, \beta) = \mathcal{N}(\vec{t} | \vec{0}, C), \text{ where}$$

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

$N \times N$

More explicitly,

$$\log p(\vec{E}|X, \vec{\alpha}, \beta) = \frac{N}{2} \log \beta + \frac{1}{2} \sum_{i=1}^M \log \alpha_i - \frac{1}{2} \vec{E}^T C^{-1} \vec{E} + \frac{1}{2} \log |\Sigma| - \frac{N}{2} \log(2\pi)$$

$$\log |C| = \log |\Sigma| - N \log \beta - \sum_{i=1}^M \log \alpha_i$$

[see next page]

Then, first of all

$$\frac{1}{2} \vec{E}^T C^{-1} \vec{E} = \frac{1}{2} \vec{E}^T (\beta^{-1} \mathbb{I} + \Phi A^{-1} \Phi^T)^{-1} \vec{E} \quad \textcircled{1}$$

Woodbury identity:

$$\left( \underbrace{A}_{\beta^{-1} \mathbb{I}} + \underbrace{B D^{-1} C}_{A^{-1} \Phi^T} \right)^{-1} = A^{-1} - A^{-1} B (D + C A^{-1} B)^{-1} C A^{-1}$$

$$\begin{aligned} \textcircled{1} \frac{1}{2} \vec{E}^T (\beta \mathbb{I} - \beta \Phi (A + \beta \Phi^T \Phi)^{-1} \Phi^T \beta) \vec{E} &= \\ &= \frac{1}{2} (\beta \vec{E}^T \vec{E} - \underbrace{\beta \vec{E}^T \Phi}_{\vec{m}^T} \underbrace{\Sigma^{-1} \Sigma \Phi^T \vec{E}}_{\vec{m}} \beta) = \\ &= \frac{1}{2} (\beta \vec{E}^T \vec{E} - \vec{m}^T \Sigma^{-1} \vec{m}) \end{aligned}$$

Next,  $\frac{\partial}{\partial \alpha_j} \log p(\vec{E}|X, \vec{\alpha}, \beta) = \frac{1}{2 \alpha_j} + \frac{1}{2} \frac{\partial}{\partial \alpha_j} \log |\Sigma| \quad \textcircled{2}$

$$\Sigma^{-1} \text{Tr} \left( \Sigma^{-1} \frac{\partial \Sigma}{\partial \alpha_j} \right) \quad \textcircled{3}$$

$$\textcircled{3} - \text{Tr} \left( \Sigma^{-1} \Sigma \frac{\partial}{\partial \alpha_j} (A + \beta \Phi^T \Phi) \Sigma \right) = - \sum_{jj} \alpha_j$$

$\mathbb{I}_{(j,j)}$ : matrix with 1 at  $(j,j)$ , 0's everywhere else

$\oplus \frac{1}{2} m^T \beta \Sigma^{-1} m = \frac{1}{2} \sum d_j + \frac{1}{2} \sum d_j + \frac{1}{2} m^2$

matrix with  $\lambda$  at  $(j,j)$ , 0's elsewhere

Then  $1 - d_j \sum d_j = -d_j m_j^2$

$d_j = \frac{d_j \sum d_j - 1}{m_j^2}$

off by a minus sign (7.87)

Double-check  $\log |C|$ :

$$\begin{aligned}
 |C| &= |\beta^{-1} \mathbb{I}_N + \Phi A^{-1} \Phi^T| = |(\beta^{-1} \mathbb{I}_N)(\mathbb{I}_N + \beta \Phi A^{-1} \Phi^T)| = \\
 &= |\beta^{-1} \mathbb{I}_N| |\mathbb{I}_N + \beta \underbrace{\Phi A^{-1} \Phi^T}_{"A" "B^T"}| = \beta^{-N} |\mathbb{I}_M + \beta A^{-1} \Phi^T \Phi| = \\
 &= \beta^{-N} |A^{-1} (A + \beta \Phi^T \Phi)| = \beta^{-N} |A^{-1}| |\Sigma^{-1}| \ominus \\
 &\quad \uparrow \quad \uparrow \quad \uparrow \\
 |AB| &= |A||B| \quad \beta^{-N} \quad \uparrow \\
 |\mathbb{I}_N + AB^T| &= |\mathbb{I}_M + A^T B| \quad \uparrow \\
 &= \beta^{-N} |A^{-1}| |\Sigma^{-1}| \ominus \\
 &\quad \uparrow \\
 &= \frac{1}{\prod_i d_i}
 \end{aligned}$$

$\ominus \frac{1}{\beta^N \prod_i d_i} \frac{1}{|\Sigma|}$ , so that

$\log |C| = -N \log \beta - \sum_{i=1}^M \log d_i - \log |\Sigma|$ ,  
 just as before

$$\textcircled{+} \quad \frac{1}{2} \frac{\partial \bar{m}^T}{\partial \alpha_j} \Sigma^{-1} \bar{m} + \frac{1}{2} \bar{m}^T \underbrace{\frac{\partial \Sigma^{-1}}{\partial \alpha_j}}_{\Pi_{(j,j)}} \bar{m} + \frac{1}{2} \bar{m}^T \Sigma^{-1} \frac{\partial \bar{m}}{\partial \alpha_j} \diamond$$

Note that

$$\frac{\partial \bar{m}}{\partial \alpha_j} = \beta \frac{\partial (\Sigma^{-1})^{-1}}{\partial \alpha_j} \Phi^T \bar{t} = -\beta \Sigma \underbrace{\frac{\partial (\Sigma^{-1})}{\partial \alpha_j}}_{\Pi_{(j,j)}} \Sigma \Phi^T \bar{t}, \text{ s.t.}$$

$$\begin{aligned} \frac{1}{2} \bar{m}^T \Sigma^{-1} \frac{\partial \bar{m}}{\partial \alpha_j} &= -\frac{1}{2} \bar{m}^T \Sigma^{-1} \Sigma \Pi_{(j,j)} \underbrace{(\beta \Sigma \Phi^T \bar{t})}_{\bar{m}} \\ &= -\frac{1}{2} \bar{m}^T \Pi_{(j,j)} \bar{m} = -\frac{1}{2} m_j^2. \end{aligned}$$

$$\diamond \quad \frac{1}{2 \alpha_j} - \frac{1}{2} \Sigma_{jj} \bullet - \frac{1}{2} m_j^2 + \frac{1}{2} m_j^2 - \frac{1}{2} m_j^2 =$$

$$= \frac{1}{2 \alpha_j} - \frac{1}{2} \Sigma_{jj} - \frac{1}{2} m_j^2 = 0.$$

Then  $1 - \alpha_j \Sigma_{jj} = \alpha_j m_j^2$ , or

$$(*) \quad \alpha_j = \frac{\overbrace{1 - \alpha_j \Sigma_{jj}}^{\approx \delta_j}}{m_j^2} \Leftarrow (7.87)$$

Use (\*) as an update eq'n:

$$\alpha_j^{\text{new}} = \frac{\delta_j}{\underbrace{m_j^2}_{\text{evaluated at } \alpha^{\text{old}}, \beta^{\text{old}}}}$$

Similarly,

$$\frac{\partial}{\partial \beta} \log p(\vec{t} | X, \vec{\alpha}, \beta) = 0 \quad \text{yields} \quad \text{(proof left as an exercise)}$$

$$\frac{1}{\beta_{\text{new}}} = \frac{\|\vec{t} - \Phi \vec{m}\|^2}{N - \sum_{i=1}^M \gamma_i} \leftarrow (7.88)$$

evaluated at  $\vec{\alpha}^{\text{old}}, \beta^{\text{old}}$

Algorithm: choose initial  $(\vec{\alpha}, \beta)$

⇓  
evaluate  $(\vec{m}, \Sigma)$

⇓  
update  $(\vec{\alpha}, \beta)$

[Stop upon convergence,  
record  $(\vec{\alpha}, \beta)$ ]

It turns out that many  $\alpha_i$ 's are driven to very large values, so that posteriors for the corresponding  $w_i$ 's are peaked at  $\emptyset$  (with variance  $\emptyset$ ), and thus the  $\gamma_i(\vec{x})$  are removed from the model.

If we use  $y(\vec{x}_\bullet) = \sum_{n=1}^N w_n k(\vec{x}, \vec{x}_n) + b,$

the datapoints  $\vec{x}_n$  for which  $w_n \neq 0$  are called relevance vectors.

Thus sparse models are created automatically.

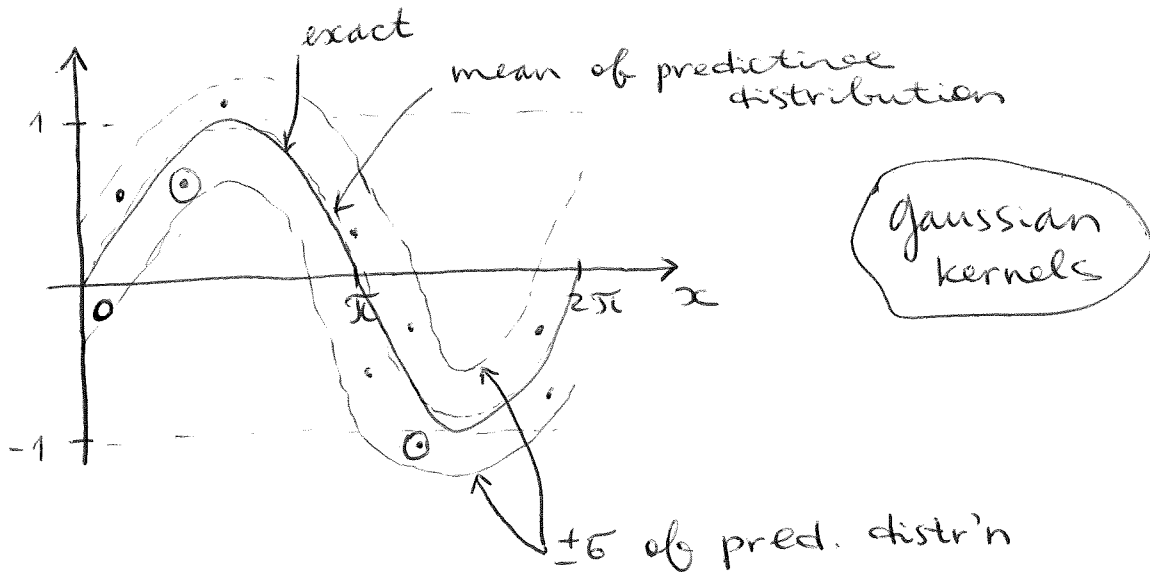
Finally,

$$p(t | \vec{x}, X, \vec{t}, \vec{\alpha}, \hat{\beta}) = \int d\vec{w} p(t | \vec{x}, \vec{w}, \hat{\beta}) p(\vec{w} | X, \vec{t}, \vec{\alpha}, \hat{\beta}) \quad \textcircled{1}$$

$$\textcircled{2} \mathcal{N}(t | \vec{m}^T \vec{\gamma}(\vec{x}), \sigma^2(\vec{x})), \text{ where}$$

$$\sigma^2(\vec{x}) = \frac{1}{\beta} + \vec{\gamma}^T(\vec{x}) \Sigma \vec{\gamma}(\vec{x}) \leftarrow \text{analogous to (3.58), (3.59)}$$

Ex:



⊙ are relevance vectors, there are relatively few.

### Sparsity analysis

Ex: consider  $N=2$ , s.t. we have  $\vec{t} = \overline{t_1 t_2}$ ,  
and a single basis  $f^n$   $\Psi(\vec{x}) \Rightarrow$  single  $d$ .  
Measurement noise is described by  $\beta$ .

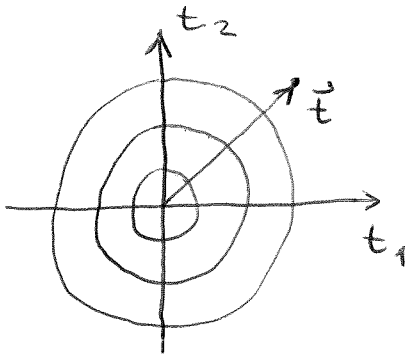
Then  $p(\vec{t} | \alpha, \beta) = \mathcal{N}(\vec{t} | \vec{0}, C)$ , where

$$\underbrace{C}_{2 \times 2} = \beta^{-1} \mathbb{I} + \alpha^{-1} \underbrace{\vec{\Psi} \cdot \vec{\Psi}^T}_{N \times N} \Leftrightarrow \begin{matrix} \text{in general,} \\ C = \beta^{-1} \mathbb{I} + \alpha A^{-1} \alpha^T \end{matrix}$$

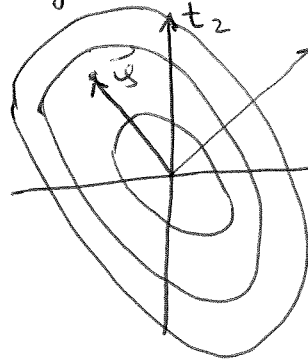
$\vec{\Psi} = \overline{\Psi(\vec{x}_1), \Psi(\vec{x}_2)}$

We maximize  $p(\vec{t} | \alpha, \beta)$  to find  $\hat{\alpha}, \hat{\beta}$ .

$$\lambda = \infty : C = \beta^{-1} \mathbb{I}$$



$\lambda$  finite but  $\vec{\psi}_i$  poorly aligned with  $\vec{t}$ :



prob. of observing  $\vec{t}$  decreases, so that  $\lambda = \infty$  is preferable

Consider  $M$  basis functions.

Idea: make dependence on a specific  $\lambda_i$  explicit in  $C$  &  $\log p(\vec{t} | X, \vec{\lambda}, \beta) \equiv L(\vec{\lambda})$ .

$$C = \beta^{-1} \mathbb{I} + \lambda_i^{-1} \vec{\psi}_i \vec{\psi}_i^T + \sum_{j \neq i} \lambda_j^{-1} \vec{\psi}_j \vec{\psi}_j^T \equiv$$

$$\equiv C_{-i} + \lambda_i^{-1} \vec{\psi}_i \vec{\psi}_i^T$$

$$\vec{\psi}_i = \begin{pmatrix} \psi_i(\vec{x}_1) & \dots & \psi_i(\vec{x}_N) \end{pmatrix}$$

$$C_{-i}^{-1} = (C_{-i}^{-1})^T$$

$$|C| = |C_{-i}| |\mathbb{I}_N + \underbrace{C_{-i}^{-1} \lambda_i^{-1} \vec{\psi}_i}_{\vec{a}} \underbrace{\vec{\psi}_i^T}_{\vec{b}^T}| = |C_{-i}| (1 + \lambda_i^{-1} \vec{\psi}_i^T C_{-i}^{-1} \vec{\psi}_i)$$

$$|\mathbb{I}_N + \underbrace{\vec{a} \vec{b}^T}_{\text{outer product}}| = 1 + \underbrace{\vec{a}^T \vec{b}}_{\text{inner product}}$$

$$C^{-1} = \left( \underbrace{C_{-i}}_A + \lambda_i^{-1} \underbrace{\vec{\psi}_i}_{B} \underbrace{\vec{\psi}_i^T}_C \right)^{-1} \stackrel{\text{Woodbury identity}}{=} C_{-i}^{-1} - C_{-i}^{-1} \vec{\psi}_i (\lambda_i + \vec{\psi}_i^T C_{-i}^{-1} \vec{\psi}_i)^{-1} \vec{\psi}_i^T C_{-i}^{-1} \quad \text{number}$$

$$\equiv C_{-i}^{-1} - \frac{C_{-i}^{-1} \vec{\psi}_i \vec{\psi}_i^T C_{-i}^{-1}}{\lambda_i + \vec{\psi}_i^T C_{-i}^{-1} \vec{\psi}_i}$$



Next,  $L(\vec{\lambda}) = -\frac{1}{2} \vec{t}^T C^{-1} \vec{t} - \frac{N}{2} \log(2\pi) - \frac{1}{2} \log |C| =$

$$= -\frac{N}{2} \log(2\pi) - \frac{1}{2} \log \left\{ |C_{-i}| (1 + \lambda_i^{-1} \vec{\Psi}_i^T C_{-i}^{-1} \vec{\Psi}_i) \right\} -$$

$$-\frac{1}{2} \vec{t}^T \left\{ C_{-i}^{-1} - \frac{C_{-i}^{-1} \vec{\Psi}_i \vec{\Psi}_i^T C_{-i}^{-1}}{\lambda_i + \vec{\Psi}_i^T C_{-i}^{-1} \vec{\Psi}_i} \right\} \vec{t} =$$

$$= -\frac{1}{2} \left\{ N \log(2\pi) + \log |C_{-i}| + \vec{t}^T C_{-i}^{-1} \vec{t} \right\} \ominus$$

"  $L(\vec{\lambda}_{-i})$ ,  $\log \mathcal{J}$  with  $\vec{\Psi}_i$  omitted

$$\ominus \frac{1}{2} \log (1 + \lambda_i^{-1} \vec{\Psi}_i^T C_{-i}^{-1} \vec{\Psi}_i) + \frac{1}{2} \vec{t}^T \frac{C_{-i}^{-1} \vec{\Psi}_i \vec{\Psi}_i^T C_{-i}^{-1}}{\lambda_i + \vec{\Psi}_i^T C_{-i}^{-1} \vec{\Psi}_i} \vec{t} =$$

$$= L(\vec{\lambda}_{-i}) + \frac{1}{2} \log \left( 1 + \frac{S_i}{\lambda_i} \right) + \frac{1}{2} \frac{q_i^2}{\lambda_i + S_i} \quad \textcircled{=}$$

$$\begin{cases} S_i = \vec{\Psi}_i^T C_{-i}^{-1} \vec{\Psi}_i, & \text{sparsity of } \vec{\Psi}_i \\ q_i = \vec{\Psi}_i^T C_{-i}^{-1} \vec{t}. & \text{quality of } \vec{\Psi}_i \end{cases}$$

$$\textcircled{=} L(\vec{\lambda}_{-i}) + \frac{1}{2} \left[ \log \lambda_i - \log(\lambda_i + S_i) + \frac{q_i^2}{\lambda_i + S_i} \right]$$

"  $\lambda(\lambda_i)$ , contains all dependence on  $\lambda_i$

Then  $\frac{\partial L(\vec{\lambda})}{\partial \lambda_i} = \frac{d\lambda(\lambda_i)}{d\lambda_i} = \frac{1}{2} \left[ \frac{1}{\lambda_i} - \frac{1}{\lambda_i + S_i} - \frac{q_i^2}{(\lambda_i + S_i)^2} \right] = \frac{1}{2} \left[ \frac{(\lambda_i + S_i)^2 - \lambda_i(\lambda_i + S_i) - q_i^2 \lambda_i}{(\lambda_i + S_i)^2 \lambda_i} \right] =$

$$= \frac{1}{2} \frac{(\lambda_i^2 + S_i^2 + 2\lambda_i S_i) - \lambda_i^2 - \lambda_i S_i - q_i^2 \lambda_i}{(\lambda_i + S_i)^2 \lambda_i} \quad \textcircled{=}$$

$$\diamond \frac{1}{2} \frac{s_i^2 + \lambda_i s_i - q_i^2 \lambda_i}{(\lambda_i + s_i)^2 \lambda_i} = \frac{\lambda_i^{-1} s_i^2 - (q_i^2 - s_i)}{2(\lambda_i + s_i)^2} = 0$$

If  $q_i^2 \leq s_i$ , only  $\lambda_i \rightarrow \infty$  is a solution.

$$\lambda_i \geq 0$$

$$\text{If } q_i^2 > s_i, \quad \underbrace{\lambda_i^{-1} s_i^2}_{>0} = \underbrace{q_i^2 - s_i}_{>0} \Rightarrow \lambda_i = \frac{s_i^2}{q_i^2 - s_i} \quad (**)$$

for finite  $\lambda_i$

So,  $q_i^2$  vs.  $s_i$  determine whether  $\vec{y}_i$  will be pruned from the model or not.

Furthermore, (\*\*) gives  $\lambda_i$  given  $\vec{y}_i$ .

This gives rise to a sequential sparse Bayesian learning algorithm:

1. Initialize  $\beta$

2. Set  $\lambda_1$  for  $\vec{y}_1$  using (\*\*); set  $\lambda_{j \neq 1}$  to  $\infty$  so that only  $\vec{y}_1$  is included in the model

Repeat until convergence:

3. Compute  $\Sigma, \vec{m}, \{q_i\}, \{s_i\}$ .

4. Randomly select  $\vec{y}_i$ .

$\rightarrow$  If  $q_i^2 > s_i$  &  $\lambda_i < \infty \Rightarrow \vec{y}_i$  already in the model  $\Rightarrow$  update  $\lambda_i$  using (\*\*)

$\rightarrow$  If  $q_i^2 > s_i$  &  $\lambda_i = \infty \Rightarrow$  add  $\vec{y}_i$  to the model, update  $\lambda_i$  using (\*\*)

→ if  $q_i^2 \leq S_i$  &  $\lambda_i < \infty \Rightarrow$  remove  $\vec{q}_i$   
from the model, set  $\lambda_i = \infty$

→ if  $q_i^2 \leq S_i$  &  $\lambda_i = \infty \Rightarrow$  do nothing

5. Update  $\beta$  using (7.88) [see top of p.6]

---

One can show that this algorithm scales  
as  $\mathcal{O}(M^3)$ , where  $M$  is the number of  
active  $\vec{q}_i$ 's, typically  $\ll N$ .