

Relevance vector machines (RVMs) } Lecture 24

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

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

Regression RVMs

Recall that in regression,

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

" " $\underbrace{\hspace{10em}}_{\epsilon^2}$

[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:

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

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

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

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

prior

separate α_i for each w_i

$$\vec{\alpha} = \underline{\alpha_1 \dots \alpha_M}$$

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

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

$$A = \begin{pmatrix} \alpha_1 & & 0 \\ & \ddots & \\ 0 & & \alpha_M \end{pmatrix}$$

For RVMS, $\Phi \Rightarrow K$, with elements $k(\vec{x}_n, \vec{x}_m)$

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

likelihood, gaussian prior, gaussian

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

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

N x N

More explicitly,

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

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

Then, first of all

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

Woodbury identity:

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

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

$$\text{Next, } \frac{\partial}{\partial \lambda_j} \log p(\vec{E}|X, \vec{x}, \beta) = \frac{1}{2 \lambda_j} + \frac{1}{2} \frac{\partial}{\partial \lambda_j} \log |\Sigma| \quad \oplus$$

$$\text{Tr} \left(\Sigma^{-1} \frac{\partial \Sigma}{\partial \lambda_j} \right) \quad \boxminus$$

$$\boxminus - \text{Tr} \left(\Sigma^{-1} \Sigma \frac{\partial}{\partial \lambda_j} (A + \beta \Phi^T \Phi) \Sigma \right) = - \sum_{jj} \mathbb{I}_{(jj)}$$

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

$\oplus \frac{1}{2} m^T \Sigma^{-1} m = \frac{1}{2} \sum_j \alpha_j + \frac{1}{2} \sum_j m_j^2$

matrix with α_j at (j,j) , 0's elsewhere

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

$\alpha_j = \frac{\alpha_j \sum_{jj} - 1}{m_j^2}$

off by a minus wrt (7.87)

Double-check $\log |C|$:

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

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

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

just as before

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

Note that

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

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

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

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

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

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

Use (*) as an update eq'n:

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

Similarly,

$$\frac{\partial}{\partial \beta} \log p(\vec{t} | X, \vec{\alpha}, \beta) = 0$$

yields

(proof left as an exercise)

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

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

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

↓
evaluate (\vec{m}, Σ)

↓
update $(\vec{\alpha}, \beta)$

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

It turns out that many α_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.

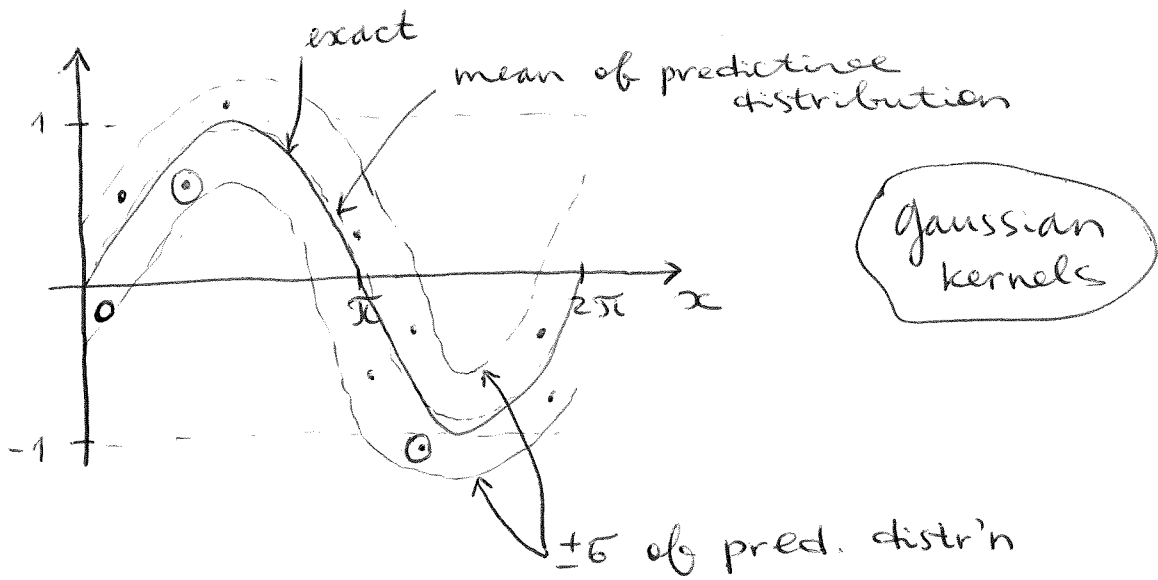
Finally,

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

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

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

Ex.:



⊙ are relevance vectors, there are relatively few.

Sparsity analysis

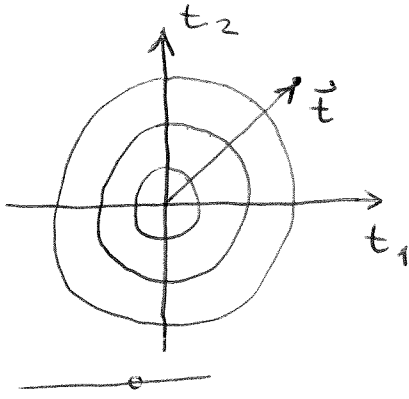
Ex.: consider $N=2$, s.t. we have $\vec{t} = \overbrace{t_1 t_2}^{(M=1)^2}$, and a single basis f'n $\mathcal{G}(\vec{x}) \Rightarrow$ single d. Measurement noise is described by β .

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

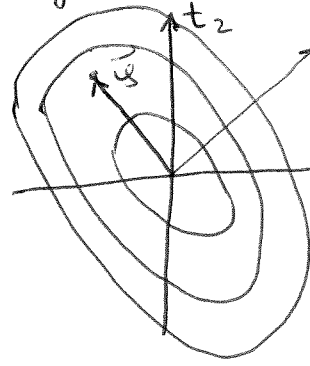
$$\underbrace{C}_{2 \times 2} = \beta^{-1} \mathbb{I} + \mathcal{L}^{-1} \underbrace{\vec{y} \cdot \vec{y}^T}_{N \times N} \Leftrightarrow \begin{matrix} \text{in general,} \\ C = \beta^{-1} \mathbb{I} + \Phi A^{-1} \Phi^T \end{matrix}$$
$$\vec{y} = \overbrace{\mathcal{G}(\vec{x}_1), \mathcal{G}(\vec{x}_2)}$$

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

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



λ finite but \vec{y}_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 λ_i explicit in C & $\log p(\vec{t} | X, \vec{\lambda}, \beta) \equiv L(\vec{\lambda})$.

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

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

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

$$C_{-i}^{-1}$$

$$\equiv (C_{-i}^{-1})^T$$

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

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

$$C^{-1} = \left(\underbrace{C_{-i}}_{\text{"A"}} + \underbrace{\lambda_i^{-1}}_{\text{"D^{-1}}} \underbrace{\vec{y}_i}_{\text{"B"}} \underbrace{\vec{y}_i^T}_{\text{"C"}} \right)^{-1} \stackrel{\text{Woodbury identity}}{=} C_{-i}^{-1} - C_{-i}^{-1} \vec{y}_i (\lambda_i + \vec{y}_i^T C_{-i}^{-1} \vec{y}_i)^{-1} \vec{y}_i^T C_{-i}^{-1} \quad \textcircled{=}$$

$$\textcircled{=} C_{-i}^{-1} - \frac{C_{-i}^{-1} \vec{y}_i \vec{y}_i^T C_{-i}^{-1}}{\lambda_i + \vec{y}_i^T C_{-i}^{-1} \vec{y}_i}$$

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

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

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

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

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

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

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

\Rightarrow

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

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

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

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

$$- \frac{q_i^2}{(\lambda_i + S_i)^2} \right] = \frac{1}{2} \left[\frac{(\lambda_i + S_i)^2 - \lambda_i(\lambda_i + S_i) - q_i^2 \lambda_i}{(\lambda_i + S_i)^2 \lambda_i} \right] =$$

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

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

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 λ_i

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

Furthermore, (**) gives λ_i given \bar{Y}_i .
This gives rise to a sequential sparse Bayesian learning algorithm:

1. Initialize β

2. Set λ_1 for \bar{Y}_1 using (**); set $\lambda_{j \neq 1}$ to ∞ so that only \bar{Y}_1 is included in the model

Repeat until convergence:

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

4. Randomly select \bar{Y}_i .

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

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

→ If $q_i^2 \leq S_i$ & $d_i < \infty \Rightarrow$ remove \vec{y}_i
from the model, set $d_i = \infty$

→ If $q_i^2 \leq S_i$ & $d_i = \infty \Rightarrow$ do nothing

5. Update β

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