

[Paper review 19]

Relevance Vector Machine Explained

(Tristan Fletcher, 2010)

[Contents]

1. Introduction
2. RVM for Regression
3. Analysis of Sparsity
4. RVM for Classification

1. Introduction

SVM

- Not a probabilistic prediction
- Only Binary decision
- have to tune hyperparameter C

RVM is more sparse, and can solve three problems above.

2. RVM for Regression

RVM = Linear Model + Modified prior for sparse solution

2.1 Model setup

1) conditional distribution : $p(t | x, w, \beta) = N(t | y(x), 1/\beta)$

2) prior :

- (LM) $p(w_i) = N(0, 1/\alpha)$
- (RVM) $p(w_i) = N(0, 1/\alpha_i)$

3) posterior : $p(w | t, X, \alpha, \beta) = N(w | m, \Sigma)$, where

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

2.2 Maximize Marginal Likelihood

- find optimal α and β by maximizing marginal likelihood, $p(t | X, \alpha, \beta)$

$$\begin{aligned} p(t | X, \alpha, \beta) &= \int p(t | X, w, \beta) p(w | \alpha) dw \\ &= \int N\left(t | w^\top \phi(x), \frac{1}{\beta}\right) N(w | 0, A^{-1}) dw \\ &= N(\mathbf{t} | 0, C) \end{aligned}$$

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

Woodbury identity

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

In our case, $A = \beta^{-1}I$, $B = \Phi$, $D = A$, $C = \Phi^T$

If we solve...

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

Solution :

$$\alpha_i^{\text{new}} = \frac{r_i}{m_i^2} = (1 - \alpha_i \Sigma_{ii}) / m_i^2 : \text{implicit}$$

$$(\beta^{\text{Nex}})^{-1} = \|t - \Phi m\|^2 / (N - \sum_i r_i)$$

- can not solve α_i^{new} directly...
- step 1) initialize α_0 and β_0
- step 2) find posterior
- step 3) update α and β
- step 4) repeat step 2 & 3

Relevance Vector

- data(vector) with non-zero weight
- $\alpha_i \approx \infty$, $w_i = 0$

RVM vs SVM

- 1) Sparsity : RVM > SVM
- 2) Generalization : RVM > SVM
- 3) Need to estimate hyperparameter : only SVM
- 4) Training Time : RVM >> SVM

3. Analysis of Sparsity

Alternative way to train RVM, due to long training time.

[Log Marginal Likelihood ($L(\alpha)$)]

$$L(\alpha) = L(\alpha_i) + \lambda(\alpha_i)$$

$$L(\alpha) = \ln(p(\mathbf{t} | \mathbf{X}, \alpha, \beta)) = \ln(N(\mathbf{t} | 0, C))$$

$$\text{where } C = \beta^{-1}I + \sum_{j \neq i} \alpha_j^{-1} \phi_j \phi_j^T + \alpha_i^{-1} \phi_i \phi_i^T = C_{-i} + \alpha_i^{-1} \phi_i \phi_i^T$$

Solution :

$$\lambda(\alpha_i) = \frac{1}{2} \left\{ \ln(|\mathbf{1} + \alpha_i^{-1} \phi_i^T C_{-i}^{-1} \phi_i|) - \mathbf{t}^T \left(\frac{C_{-i}^{-1} \phi_i \phi_i^T C_{-i}^{-1}}{\alpha_i + \phi_i^T C_{-i}^{-1} \phi_i} \right) \mathbf{t} \right\}$$

$$s_i = \phi_i^T C_{-i}^{-1} \phi_i$$

- sparsity of ϕ_i
- overlap between ϕ_i and ϕ_j

$$q_i = \phi_i^T C_{-i}^{-1} \mathbf{t}$$

- quality of ϕ_i
- $C_{-i}^{-1} \mathbf{t}$: prediction error $\rightarrow q_i$: information about ϕ_i

$$s_i > q_i \rightarrow \phi_i = 0$$

$$s_i < q_i \rightarrow \phi_i \neq 0$$

$$\frac{\partial}{\partial \alpha} L(\alpha) = \frac{\partial}{\partial \alpha} \lambda(\alpha) = \frac{\alpha_i^{-1} s_i^2 + s_i - q_i^2}{(\alpha_i + s_i)^2} 0$$

- if $s_i \leq q_i^2 \rightarrow \alpha_i = \frac{s_i^2}{s_i - q_i^2}$
- if $s_i > q_i^2 \rightarrow \alpha_i = \infty$

4. RVM for Classification

Binary case

- $y(x) = \sigma(w^T \varphi(x))$

Multi-class case

- $y_K(x) = \frac{\exp(w_K^T \varphi(x))}{\sum_j \exp(w_j^T \varphi(x))}$

Hard to calculate marginal likelihood, so use Laplace Approximation