

RVM

Атанов Андрей

HSE

3 Октября 2016

Немного про RVM

Немного про RVM

Данные

$X = \{x_i\}_{i=1}^n$, $x_i \in \mathbb{R}^d$ - объекты выборки и $t \in \mathbb{R}^n$ - вектор целевых переменных.

Немного про RVM

Данные

$X = \{x_i\}_{i=1}^n, x_i \in \mathbb{R}^d$ - объекты выборки и $t \in \mathbb{R}^n$ - вектор целевых переменных.

Посыл

$$t_n \sim \mathcal{N}(y(w, x_i), \sigma^2)$$

Немного про RVM

Данные

$X = \{\mathbf{x}_i\}_{i=1}^n$, $\mathbf{x}_i \in \mathbb{R}^d$ - объекты выборки и $\mathbf{t} \in \mathbb{R}^n$ - вектор целевых переменных.

Посыл

$$t_n \sim \mathcal{N}(y(\mathbf{w}, \mathbf{x}_i), \sigma^2)$$

$$y(\mathbf{w}, \mathbf{x}) = \sum_{i=1}^D w_i \phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x})$$

$\boldsymbol{\phi}(\mathbf{x}) = (\phi_1(\mathbf{x}), \dots, \phi_D(\mathbf{x}))$ - вектор признаков.

Немного про RVM

Функция правдоподобия

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \sigma^2) = (2\pi\sigma^2)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\Phi\mathbf{w} - \mathbf{t}\|^2 \right\} \Rightarrow \mathbf{w}_{\text{ML}}$$

Немного про RVM

Функция правдоподобия

$$p(t|X, w, \sigma^2) = (2\pi\sigma^2)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\Phi w - t\|^2 \right\} \Rightarrow w_{ML}$$

Есть серьезная проблема - переобучение.

Немного про RVM

Функция правдоподобия

$$p(t|X, w, \sigma^2) = (2\pi\sigma^2)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\Phi w - t\|^2 \right\} \Rightarrow w_{ML}$$

Есть серьезная проблема - переобучение.

Регуляризация

$$w \sim \mathcal{N}(0, A^{-1})$$

$$A = \text{diag}(\alpha_1, \dots, \alpha_D)$$

Немного про RVM

Функция правдоподобия

$$p(t|X, w, \sigma^2) = (2\pi\sigma^2)^{n/2} \exp \left\{ -\frac{1}{2\sigma^2} \|\Phi w - t\|^2 \right\} \Rightarrow w_{ML}$$

Есть серьезная проблема - переобучение.

Регуляризация

$$w \sim \mathcal{N}(0, A^{-1})$$

$$A = \text{diag}(\alpha_1, \dots, \alpha_D)$$

Апостериорное распределение

$$p(w|t) \propto p(t|w)p(w|\alpha) \Rightarrow w_{MAP}$$

Немного про RVM

Хотим

Зная обучающую выборку $\{X, t\}$, оценить для нового объекта x^* целевую переменную t^* (восстановить плотность).

$$p(t^*|t) = \int p(t^*|w, \alpha, \sigma^2)p(\alpha, \sigma^2, w|t)dw d\alpha d\sigma^2$$

Немного про RVM

Хотим

Зная обучающую выборку $\{X, t\}$, оценить для нового объекта x^* целевую переменную t^* (восстановить плотность).

$$p(t^*|t) = \int p(t^*|w, \alpha, \sigma^2)p(\alpha, \sigma^2, w|t)dw d\alpha d\sigma^2$$

Проблема

$$p(w, \alpha, \sigma^2|t) = \frac{p(t|\sigma^2, \alpha, w)p(\alpha, \sigma^2, w)}{p(t)}$$

Немного про RVM

Хотим

Зная обучающую выборку $\{X, t\}$, оценить для нового объекта x^* целевую переменную t^* (восстановить плотность).

$$p(t^*|t) = \int p(t^*|w, \alpha, \sigma^2)p(\alpha, \sigma^2, w|t)dw d\alpha d\sigma^2$$

Проблема

$$p(w, \alpha, \sigma^2|t) = \frac{p(t|\sigma^2, \alpha, w)p(\alpha, \sigma^2, w)}{p(t)}$$

$p(t)$ - нельзя вычислить аналитически.

Немного про RVM

Хотим

Зная обучающую выборку $\{X, t\}$, оценить для нового объекта x^* целевую переменную t^* (восстановить плотность).

$$p(t^*|t) = \int p(t^*|w, \alpha, \sigma^2)p(\alpha, \sigma^2, w|t)dw d\alpha d\sigma^2$$

Проблема

$$p(w, \alpha, \sigma^2|t) = \frac{p(t|\sigma^2, \alpha, w)p(\alpha, \sigma^2, w)}{p(t)}$$

$p(t)$ - нельзя вычислить аналитически.

Решение

$$p(w, \alpha, \sigma^2|t) = p(w|t, \alpha, \sigma^2)p(\alpha, \sigma^2|t)$$

Выбор модели (настройка гиперпараметров)

Максимизация обоснованности

$$p(\alpha, \sigma^2 | t) \propto p(t | \alpha, \sigma^2) p(\alpha) p(\sigma^2)$$

$$\alpha_{\text{MP}}, \sigma_{\text{MP}}^2 = \operatorname{argmax}\{p(t | \alpha, \sigma^2)\}$$

Для задачи регрессии обоснованность можно посчитать аналитически, для классификации приходится пользоваться аппроксимацией Лапласа.

Выбор модели (настройка гиперпараметров)

Максимизация обоснованности

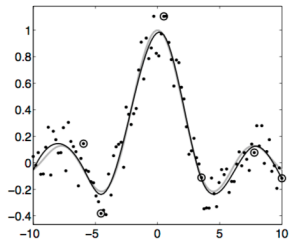
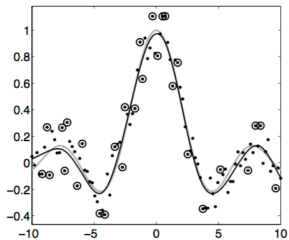
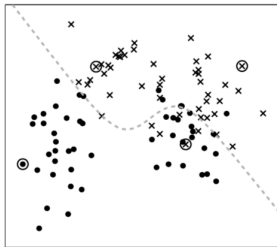
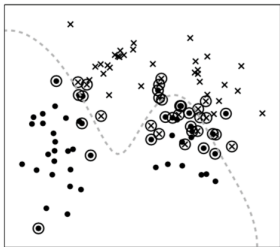
$$p(\alpha, \sigma^2 | t) \propto p(t | \alpha, \sigma^2) p(\alpha) p(\sigma^2) \\ \alpha_{\text{MP}}, \sigma_{\text{MP}}^2 = \operatorname{argmax}\{p(t | \alpha, \sigma^2)\}$$

Для задачи регрессии обоснованность можно посчитать аналитически, для классификации приходится пользоваться аппроксимацией Лапласа.

Итоговое распределение

$$p(t^* | t, \alpha_{\text{MP}}, \sigma_{\text{MP}}^2) = \int p(t^* | w, \alpha_{\text{MP}}, \sigma_{\text{MP}}^2) p(w | t, \alpha_{\text{MP}}, \sigma_{\text{MP}}^2) dw \\ = \mathcal{N}(t^* | y_*, \sigma_*^2)$$

Пример



RVM

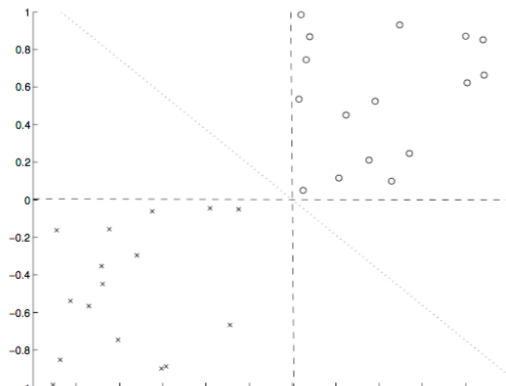
Обсуждение

- ▶ $\alpha_i \rightarrow +\infty \Rightarrow$ вес i -го признака обнуляется и соответствующий признак удаляется из модели. Обученная модель оказывается намного разреженнее по сравнению с SVM. Данный эффект носит название Auto Relevance Determination.

RVM

Обсуждение

- ▶ $\alpha_i \rightarrow +\infty \Rightarrow$ вес i -го признака обнуляется и соответствующий признак удаляется из модели.
Обученная модель оказывается намного разреженнее по сравнению с SVM. Данный эффект носит название Auto Relevance Determination.



RVM

Обсуждение

- ▶ Автоматическая настройка параметров α, σ^2 .

RVM

Обсуждение

- ▶ Автоматическая настройка параметров α, σ^2 .
- ▶ Использование всей выборке для настройки параметров и гиперпараметров

Обсуждение

- ▶ Автоматическая настройка параметров α, σ^2 .
- ▶ Использование всей выборке для настройки параметров и гиперпараметров
- ▶ НО все-таки необходимо подбирать параметры ядровой функции

Подбор параметра ядровой функции

Идея

При использовании ядровой ф-ии в качестве признаков, правдоподобие так же зависит и от параметра, входящего в нее. Например, для RBF ядра $K(\mathbf{x}_i, \mathbf{x}_j) = \exp\{-\gamma\|\mathbf{x}_i - \mathbf{x}_j\|^2\}$:

$$p(t|\alpha, \sigma^2, \gamma) \rightarrow \max_{\gamma \in C}$$

C - конечное множество.

Результаты

Кросс-валидация

Data set	Polynomial Kernel			RBF Kernel		
	CV R^2	Test R^2	$\log p(t \alpha, \sigma^2)$	CV R^2	Test R^2	$\log p(t \alpha, \sigma^2)$
Airfoil	0.509	0.591	-221	0.509	0.645	-198
Concrete	0.61	0.733	-169	0.672	0.778	-178
CCPP	0.927	0.926	-61	0.933	0.935	-42
Life expectancy	0.882	0.345	-43	0.922	0.814	-23
Friedman #3	0.482	0.734	-282	0.765	0.879	-196
Boston Housing	0.834	0.906	-308	0.862	0.919	-133

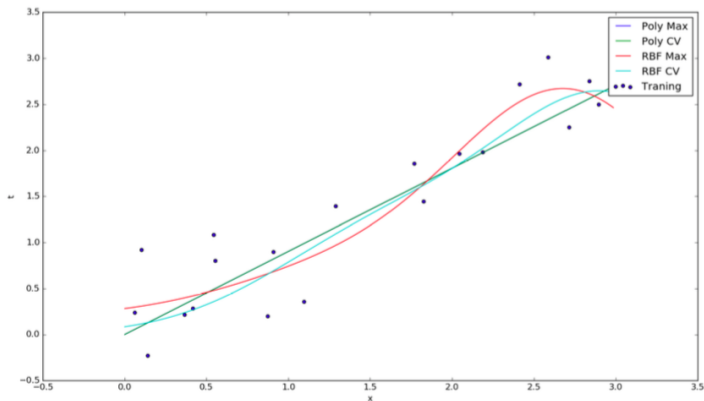
Максимизация обоснованности

Data set	Polynomial Kernel			RBF Kernel		
	CV R^2	Test R^2	$\log p(t \alpha, \sigma^2)$	CV R^2	Test R^2	$\log p(t \alpha, \sigma^2)$
Airfoil	0.495	0.591	-221	0.666	0.739	-168
Concrete	0.656	0.722	-152	0.663	0.762	-155
Poly data	0.964	0.985	21	0.969	0.991	18
CCPP	0.927	0.926	-58	0.933	0.935	-35
Life expectancy	0.724	0.328	-41	0.916	0.888	-21
Friedman #3	0.478	0.721	-269	0.049	0.095	-136
Boston Housing	0.789	0.925	-255	0.478	0.713	-40

Подходит ли ядро для данных?

Рассмотрим простой искусственный пример.

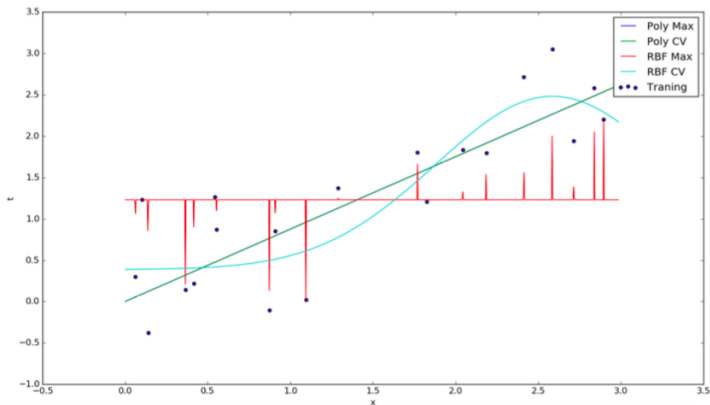
$$t = x + \mathcal{N}(0, \sigma^2)$$

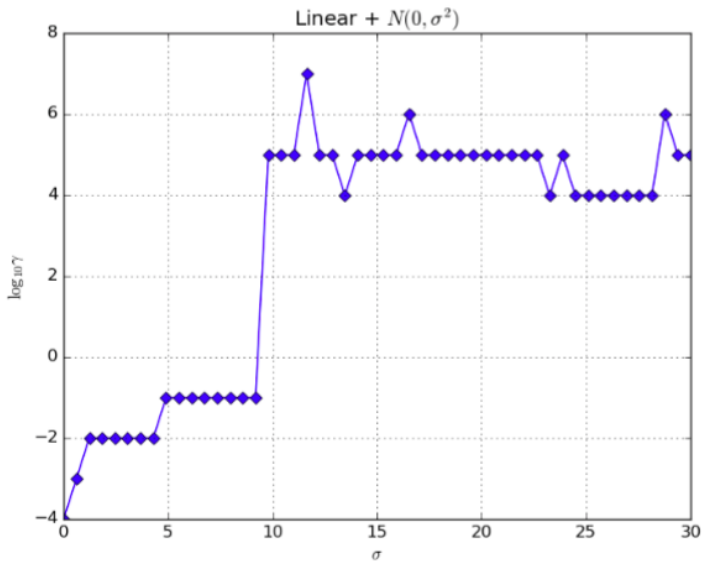


Подходит ли ядро для данных?

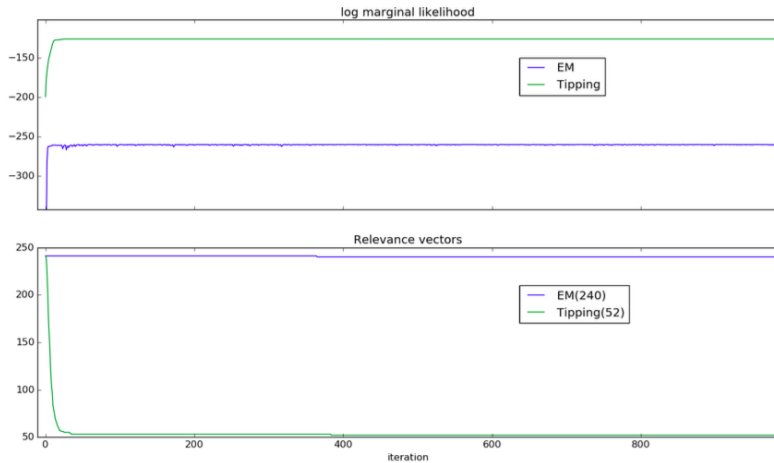
Рассмотрим простой искусственный пример.

$$t = x + \mathcal{N}(0, \sigma^2)$$

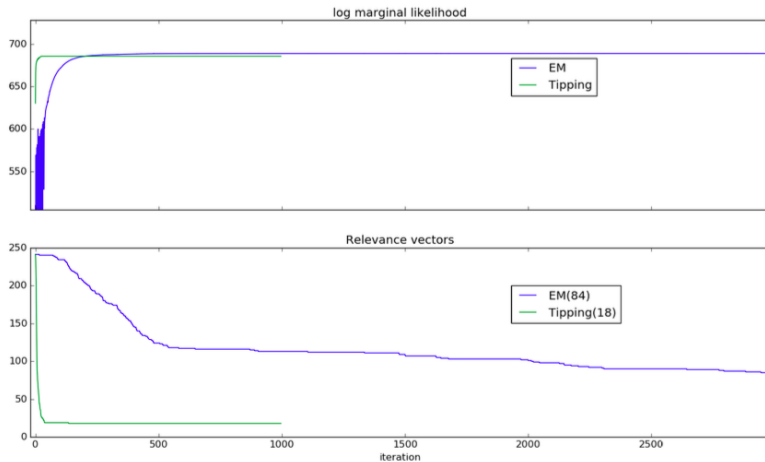




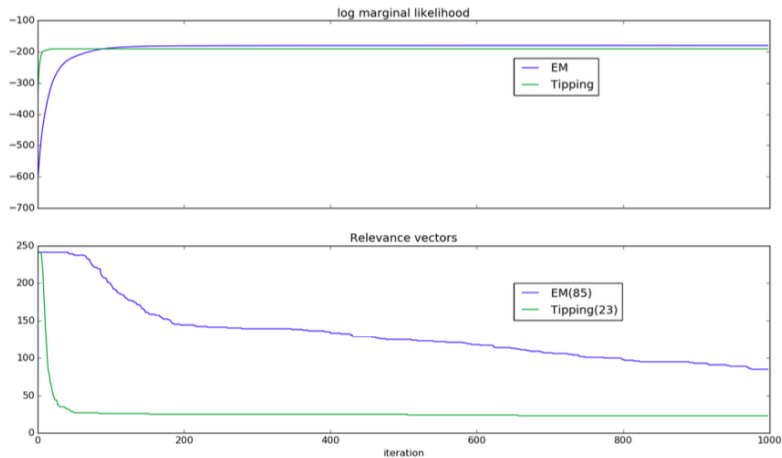
Сходимость



Сходимость



Сходимость



Литература

- ▶ Tipping, <http://www.jmlr.org/papers/volume1/tipping01a/tipping01a.pdf>.
- ▶ Bishop - 1.2-1.3, 3.3-3.4, 7.2.
- ▶ Лекции - <https://www.youtube.com/watch?v=sZxE-BrSMAE&list=PLlb7e2G7aSpR8mbaShVBods-hGaFGifkl>

Спасибо за внимание!