

# Отбор информативных признаков в методе опорных векторов

Андрющенко Анастасия Михайловна, гр. 522

Санкт-Петербургский государственный университет  
Математико-механический факультет  
Кафедра статистического моделирования

Научный руководитель: Коробейников А.И.  
Рецензент: к.ф.-м.н., д. Алексеева Н.П.



Санкт-Петербург  
2010г.

Задача отбора информативных признаков в задачах классификации способствует:

- уменьшению ошибки предсказания;
- определению значимых признаков;
- уменьшению размерности данных.

В работе используется модификация метода опорных векторов (support vector machines, SVM (Vapnik, 1995)).

# SVM. Задача классификации

Данные:

	$X_1$	$\dots$	$X_m$
$\mathbf{x}_1$	$x_1^{(1)}$	$\dots$	$x_1^{(m)}$
$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	$x_n^{(1)}$	$\dots$	$x_n^{(m)}$

$\mathbf{x}_i \in \mathbb{R}^m$  — наблюдения

$(X_1, \dots, X_m)$  — признаки

Метки:  $\mathbf{x}_i \rightarrow y_i \in \{\pm 1\}$

Задача классификации:

построить  $f : \quad f(\mathbf{x}_i) = y_i \quad \forall i = 1 \dots m.$

Задача отбора признаков:

$(X_1, \dots, X_m) \rightarrow (X_{i_1}, \dots, X_{i_l})$

# SVM. Задача классификации

Данные:

	$X_1$	$\dots$	$X_m$
$\mathbf{x}_1$	$x_1^{(1)}$	$\dots$	$x_1^{(m)}$
$\vdots$	$\vdots$		$\vdots$
$\mathbf{x}_n$	$x_n^{(1)}$	$\dots$	$x_n^{(m)}$

$\mathbf{x}_i \in \mathbb{R}^m$  — наблюдения

$(X_1, \dots, X_m)$  — признаки

Метки:  $\mathbf{x}_i \rightarrow y_i \in \{\pm 1\}$

Задача классификации:

построить  $f$  :  $f(\mathbf{x}_i) = y_i \quad \forall i = 1 \dots m.$

Задача отбора признаков:

$$(X_1, \dots, X_m) \rightarrow (X_{i_1}, \dots, X_{i_l})$$

# SVM. Задача классификации

Данные:

	$X_1$	$\dots$	$X_m$	
$\mathbf{x}_1$	$x_1^{(1)}$	$\dots$	$x_1^{(m)}$	$\mathbf{x}_i \in \mathbb{R}^m$ — наблюдения $(X_1, \dots, X_m)$ — признаки
$\vdots$	$\vdots$		$\vdots$	
$\mathbf{x}_n$	$x_n^{(1)}$	$\dots$	$x_n^{(m)}$	

Метки:  $\mathbf{x}_i \rightarrow y_i \in \{\pm 1\}$

Задача классификации:

построить  $f$  :  $f(\mathbf{x}_i) = y_i \quad \forall i = 1 \dots m.$

Задача отбора признаков:

$$(X_1, \dots, X_m) \rightarrow (X_{i_1}, \dots, X_{i_l})$$

Целевая функция SVM:

$$f_{\sigma}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \mathbf{x} \rangle - b).$$

Параметры  $(\mathbf{w}, b)$  получены решением:

$$\frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i \longrightarrow \min_{\mathbf{w}, b, \xi}$$

$$y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle - b) \geq 1 - \xi_i,$$

$$\xi_i \geq 0, \quad i = 1, \dots, n.$$

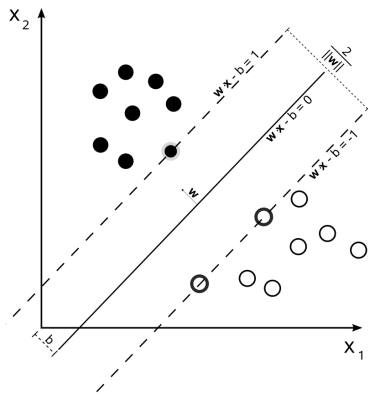
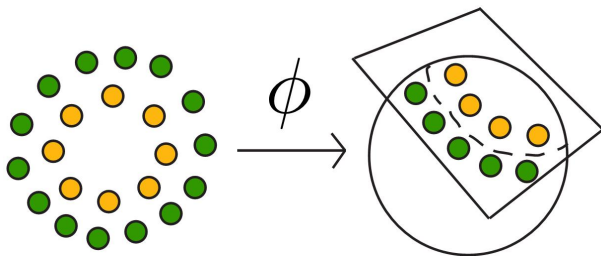


Рис.: Гиперплоскость.

Данные не всегда отделимы в исходном пространстве  
 $\Rightarrow$  они рассматриваются в спрямляющем:

$$\phi : \mathbb{R}^m \rightarrow \mathcal{H}$$



Введем ядро:

$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle_{\mathcal{H}}.$$

Целевая функция SVM:

$$f_{\sigma}(\mathbf{x}) = \text{sign}(\langle \mathbf{w}, \phi_{\sigma}(\mathbf{x}) \rangle_{\mathcal{H}} - b) = \text{sign}\left(\sum_{i=1}^{n_s} \alpha_i y_i K_{\sigma}(\mathbf{s}_i, \mathbf{x}) - b\right).$$

Примеры ядер:

- Линейное:  $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$ .
- Полиномиальное:  $K_{\gamma, r, d}(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \langle \mathbf{x}_i, \mathbf{x}_j \rangle + r)^d$ ,  $\gamma > 0$ .
- Экспоненциально-радиальное (RBF — radial basis function):  
 $K_{\gamma}(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma \|\mathbf{x}_i - \mathbf{x}_j\|^2}.$



Рассматриваются ядра:

$$K(\mathbf{x}, \mathbf{z}) = \ell \left( \|\boldsymbol{\sigma}(\mathbf{x} - \mathbf{z})\|^2 \right) = \ell \left( \sum_{k=1}^m \sigma_k^2 (x_k - z_k)^2 \right),$$

где  $\ell$  — монотонная функция. Масштабирующие коэффициенты  $\{\sigma_k\}_{k=1}^m$  отражают степень влияния признака на результат классификации.

Задача:

- 1 Реализация алгоритма пересчета масштабирующих коэффициентов (Grandvalet, Canu, 2003) на R.
- 2 Проверка и исследование на модельных данных.
- 3 Анализ реальных данных из кардиологии.

Рассматриваются ядра:

$$K(\mathbf{x}, \mathbf{z}) = \ell \left( \|\boldsymbol{\sigma}(\mathbf{x} - \mathbf{z})\|^2 \right) = \ell \left( \sum_{k=1}^m \sigma_k^2 (x_k - z_k)^2 \right),$$

где  $\ell$  — монотонная функция. Масштабирующие коэффициенты  $\{\sigma_k\}_{k=1}^m$  отражают степень влияния признака на результат классификации.

Задача:

- 1 Реализация алгоритма пересчета масштабирующих коэффициентов (Grandvalet, Canu, 2003) на R.
- 2 Проверка и исследование на модельных данных.
- 3 Анализ реальных данных из кардиологии.

Параметры классификатора:  $C, \sigma_0$ .

Обучающий критерий:

$$\begin{aligned} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^n \xi_i &\longrightarrow \min_{\sigma, \mathbf{w}, b, \xi}, \\ y_i(\langle \mathbf{w}, \phi_{\sigma}(\mathbf{x}_i) \rangle_{\mathcal{H}} - b) &\geq 1 - \xi_i, \\ \xi_i &\geq 0, \quad i = 1, \dots, n, \\ \frac{1}{m} \sum_{k=1}^m \sigma_k^2 &= \sigma_0^2. \end{aligned}$$

Вместо сложной задачи итерационно решается несколько более простых:

- 1 Зафиксировать  $\sigma$  и построить SVM.
- 2 Метод позволяет получить  $w$ ,  $\sum_{i=1}^n \xi_i$  и  $\partial \sum_{i=1}^n \xi_i / \partial \sigma$  в виде функций от  $\sigma$ :

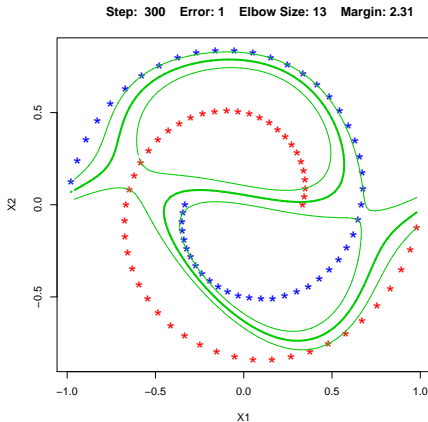
$$g(\sigma) = \frac{1}{2} \|w\|^2 + C \sum_{i=1}^n \xi_i.$$

$$g(\sigma) \rightarrow \min_{\sigma}.$$

- 3 На шаге  $l$ , начиная с  $\sigma^{(l)}$ , вычислить оптимальные  $(\hat{w}(\sigma^{(l)}), \hat{b}(\sigma^{(l)}))$ .  $\sigma^{(l+1)}$  определяется с помощью метода сопряженных градиентов.

Данные:

- 1 Координаты двух спиралей. Значимость признаков одинакова.



Данные:

- ❶ Координаты двух спиралей. Значимость признаков одинакова.
- ❷ Данные из (Chapelle, Vapnik, 2002).
  - Линейно отделимые данные. 6 признаков из 10 значимы, остальные шум.
  - Линейно не отделимые данные. 2 признака из 10 значимы, остальные шум.

Полученные масштабирующие коэффициенты  $\sigma$  соответствуют характеру данных.

Данные о 422 пациентах, перенесших операцию на открытом сердце.

Значимые признаки определялись для следующих задач:

- бинарная классификация: пациенты с наличием или отсутствием ПКТС;
- тернарная классификация: ПКТС отсутствует/ранний/поздний.

50 признаков: 22 количественных, 28 категориальных.

	Возр.	ИМТ	Кардиопл.	ФВ	Эозин.	Л1	Л7
1	67	23	13.50	68	0	8.70	7.10
2	57	20	6.30	68	1	9.90	6.70
3	54	37	9.50	64	1	14.70	7.00
4	57	28	13.50	69	0	17.20	10.40
...							

## Бинарная задача

Всего 34 значимых признака:

- 14 количественных;
- 17 категориальных;
- 3 градации у двух категориальных признаков.

Адреналин	0.04*	Дренир.1	0.21	Тип опер.2	0.11
Возраст	0.22	Дренир.2	0.13	Тип опер.3	0.07*
Аллергия	0.16	Дренир.3	0.09*	Тип опер.4	0.02*
Анемия	0.01*	Дренир.пот.	0.19	Тип опер.5	0.01*
Антиког.тер.	0.08*	Дренир.вр.	0.05*	Тип опер.6	0.06*
Переж.аорты	0.14	Фр.выброса	0.17	Веноз.застой	0.13
Арт.шунт	0.05*	Эозинофилы	0.24	Темп.прайм.	0.19
Аутоимунные	0.01*	СОЭ	0.19	МНО	0.06*
ИМТ	0.24	EuroSCORE	0.19	Апп.иск.кр.	0.14
Число шунтов	0.09*	EuroSCORE2	0.09*	Реперфузия	0.07*
Кардиоплегия	0.08*	ССН	0.14	Дыхат.недост.	0.13
Коронарогр.	0.28	Гиперглик.	0.16	Наруш.ритма	0.11

...

Значимые признаки, \* не значимые.



## Тернарная задача

Рассматривались 3 бинарные задачи. Для каждой из них получены масштабирующие коэффициенты.

Признак	n-el	l-ne	e-nl	Признак	n-el	l-ne	e-nl
Адреналин	0.01*	0.68	0.78	СОЭ	0.01*	0.01*	0.51
Возраст	0.01*	0.99	0.98	EuroSCORE	1.77	0.34	0.37
Переж.аорты	0.01*	0.25	0.01*	EuroSCORE2	1.15	0.45	0.61
ИМТ	3.16	1.02	1.18	Ин.поддержка	0.01*	0.01*	0.07*
Число шунтов	0.01*	0.31	0.19	Лейк.д.1	1.05	0.65	0.68
Кардиоплегия	0.01*	0.16	0.25	Лейк.д.7	2.35	0.63	0.18
Хрон.серд.нед.	0.01*	0.04	0.21	Мезатон	0.01*	0.01*	0.01*
Дрен.потери	0.01*	0.94	0.97	Темп.прайм.	0.01*	0.49	0.32
Дренир.вр.	3.31	0.41	0.69	МНО	1.29	0.18	0.43
Фр.выброса	4.64	0.97	0.71	Апп.иск.кр.	3.41	0.74	1.01
Эозинофилы	0.01*	0.34	0.49	Реперфузия	0.01*	0.27	0.46

Значимые признаки, \* не значимые.

## Итоги:

- Реализован алгоритм отбора информативных признаков на языке R.
- Его работа проверена и исследована на модельных данных.
- Осуществлен анализ реальных данных.

## Перспективы:

- Создание приложения.
- Комбинирование алгоритма с другими методами.