

# Задача обучения линейного классификатора

**Дано:**

Обучающая выборка  $X^\ell = (x_i, y_i)_{i=1}^\ell$ ,

$x_i$  — объекты, векторы из множества  $X = \mathbb{R}^n$ ,

$y_i$  — метки классов, элементы множества  $Y = \{-1, +1\}$ .

**Найти:**

Параметры  $w \in \mathbb{R}^n$ ,  $w_0 \in \mathbb{R}$  линейной модели классификации

$$a(x; w, w_0) = \text{sign}(\langle x, w \rangle - w_0).$$

**Критерий** — минимизация эмпирического риска:

$$\sum_{i=1}^{\ell} [a(x_i; w, w_0) \neq y_i] = \sum_{i=1}^{\ell} [M_i(w, w_0) < 0] \rightarrow \min_{w, w_0}.$$

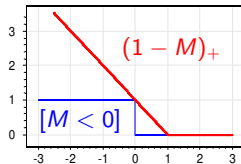
где  $M_i(w, w_0) = (\langle x_i, w \rangle - w_0) y_i$  — отступ (margin) объекта  $x_i$ ,  
 $b(x) = \langle x, w \rangle - w_0$  — дискриминантная функция.

# Аппроксимация и регуляризация эмпирического риска

Эмпирический риск — это кусочно-постоянная функция.  
Заменяем его оценкой сверху, непрерывной по параметрам:

$$\begin{aligned} Q(w, w_0) &= \sum_{i=1}^{\ell} [M_i(w, w_0) < 0] \leq \\ &\leq \sum_{i=1}^{\ell} (1 - M_i(w, w_0))_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}. \end{aligned}$$

- Аппроксимация штрафует объекты за приближение к границе классов, увеличивая зазор между классами
- Регуляризация штрафует неустойчивые решения в случае мультиколлинеарности



# Оптимальная разделяющая гиперплоскость

Линейный классификатор:

$$a(x, w) = \text{sign}(\langle w, x \rangle - w_0), \quad w, x \in \mathbb{R}^n, \quad w_0 \in \mathbb{R}.$$

Пусть выборка  $X^\ell = (x_i, y_i)_{i=1}^\ell$  линейно разделима:

$$\exists w, w_0 : \quad M_i(w, w_0) = y_i(\langle w, x_i \rangle - w_0) > 0, \quad i = 1, \dots, \ell$$

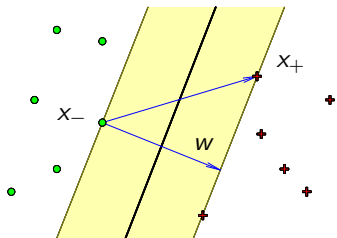
Нормировка:  $\min_{i=1, \dots, \ell} M_i(w, w_0) = 1$ .

Разделяющая полоса:

$$\{x : -1 \leq \langle w, x \rangle - w_0 \leq 1\}.$$

Ширина полосы:

$$\frac{\langle x_+ - x_-, w \rangle}{\|w\|} = \frac{2}{\|w\|} \rightarrow \max.$$



Постановка задачи в случае линейно разделимой выборки:

$$\begin{cases} \frac{1}{2} \|w\|^2 \rightarrow \min_{w, w_0}; \\ M_i(w, w_0) \geq 1, \quad i = 1, \dots, \ell. \end{cases}$$

Общий случай — линейно неразделимая выборка:

$$\begin{cases} \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \rightarrow \min_{w, w_0, \xi}; \\ M_i(w, w_0) \geq 1 - \xi_i, \quad i = 1, \dots, \ell; \\ \xi_i \geq 0, \quad i = 1, \dots, \ell. \end{cases}$$

Исключая  $\xi_i$ , получаем задачу безусловной минимизации:

$$C \sum_{i=1}^{\ell} (1 - M_i(w, w_0))_+ + \frac{1}{2} \|w\|^2 \rightarrow \min_{w, w_0}.$$

Задача математического программирования:

$$\begin{cases} f(x) \rightarrow \min_x; \\ g_i(x) \leq 0, & i = 1, \dots, m; \\ h_j(x) = 0, & j = 1, \dots, k. \end{cases}$$

Необходимые условия. Если  $x$  — точка локального минимума, то существуют множители  $\mu_i, i = 1, \dots, m, \lambda_j, j = 1, \dots, k$ :

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial x} = 0, & \mathcal{L}(x; \mu, \lambda) = f(x) + \sum_{i=1}^m \mu_i g_i(x) + \sum_{j=1}^k \lambda_j h_j(x); \\ g_i(x) \leq 0; \ h_j(x) = 0; & \text{(исходные ограничения)} \\ \mu_i \geq 0; & \text{(двойственные ограничения)} \\ \mu_i g_i(x) = 0; & \text{(условие дополняющей нежёсткости)} \end{cases}$$

Функция Лагранжа:

$$\begin{aligned}\mathcal{L}(w, w_0, \xi; \lambda, \eta) = \\ = \frac{1}{2} \|w\|^2 - \sum_{i=1}^{\ell} \lambda_i (M_i(w, w_0) - 1) - \sum_{i=1}^{\ell} \xi_i (\lambda_i + \eta_i - C),\end{aligned}$$

$\lambda_i$  — переменные, двойственные к ограничениям  $M_i \geq 1 - \xi_i$ ;

$\eta_i$  — переменные, двойственные к ограничениям  $\xi_i \geq 0$ .

$$\begin{cases} \frac{\partial \mathcal{L}}{\partial w} = 0, & \frac{\partial \mathcal{L}}{\partial w_0} = 0, & \frac{\partial \mathcal{L}}{\partial \xi} = 0; \\ \xi_i \geq 0, & \lambda_i \geq 0, & \eta_i \geq 0, \quad i = 1, \dots, \ell; \\ \lambda_i = 0 \text{ либо } M_i(w, w_0) = 1 - \xi_i, & i = 1, \dots, \ell; \\ \eta_i = 0 \text{ либо } \xi_i = 0, & i = 1, \dots, \ell; \end{cases}$$

## Двойственная задача

$$\begin{cases} -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Решив эту задачу численно относительно  $\lambda_i$ ,  
получаем линейный классификатор:

$$a(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^{\ell} \lambda_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle - w_0 \right),$$

где  $w_0 = \sum_{i=1}^{\ell} \lambda_i y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle - y_j$  для такого  $j$ , что  $\lambda_j > 0$ ,  $M_j = 1$

### Определение

Объект  $\mathbf{x}_i$  называется *опорным*, если  $\lambda_i \neq 0$ .

## Преимущества:

- Задача выпуклого квадратичного программирования имеет единственное решение.
- Выделяется множество опорных объектов.
- Имеются эффективные численные методы для SVM.
- Изящное обобщение на нелинейные классификаторы.

## Недостатки:

- Опорными объектами могут становиться выбросы.
- Нет отбора признаков в исходном пространстве  $X$ .
- Приходится подбирать константу  $C$ .



## Двойственная задача

$$\begin{cases} -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Решив эту задачу численно относительно  $\lambda_i$ ,  
получаем линейный классификатор:

$$a(\mathbf{x}) = \text{sign} \left( \sum_{i=1}^{\ell} \lambda_i y_i \langle \mathbf{x}_i, \mathbf{x} \rangle - w_0 \right),$$

где  $w_0 = \sum_{i=1}^{\ell} \lambda_i y_i \langle \mathbf{x}_i, \mathbf{x}_j \rangle - y_j$  для такого  $j$ , что  $\lambda_j > 0$ ,  $M_j = 1$

### Определение

Объект  $\mathbf{x}_i$  называется *опорным*, если  $\lambda_i \neq 0$ .

## Двойственная задача: нелинейное обобщение SVM

$$\begin{cases} -\sum_{i=1}^{\ell} \lambda_i + \frac{1}{2} \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} \lambda_i \lambda_j y_i y_j K(x_i, x_j) \rightarrow \min_{\lambda}; \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell; \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{cases}$$

Решив эту задачу численно относительно  $\lambda_i$ ,  
получаем линейный классификатор:

$$a(x) = \text{sign}\left(\sum_{i=1}^{\ell} \lambda_i y_i K(x_i, x) - w_0\right),$$

где  $w_0 = \sum_{i=1}^{\ell} \lambda_i y_i K(x_i, x_j) - y_j$  для такого  $j$ , что  $\lambda_j > 0$ ,  $M_j = 1$

### Определение

Объект  $x_i$  называется *опорным*, если  $\lambda_i \neq 0$ .

## Определение

Функция от пары объектов  $K(x, x')$  называется *ядром*, если она представима в виде скалярного произведения

$$K(x, x') = \langle \psi(x), \psi(x') \rangle$$

при некотором преобразовании  $\psi: X \rightarrow H$  из пространства признаков  $X$  в новое *спрямляющее* пространство  $H$ .

## Возможная интерпретация:

признак  $f_i(x) = K(x_i, x)$  — это оценка близости объекта  $x$  к опорному объекту  $x_i$ . Выбирая опорные объектов, SVM осуществляет отбор признаков в линейном классификаторе

$$a(x) = \text{sign} \left( \sum_{i=1}^{\ell} \lambda_i y_i K(x_i, x) - w_0 \right).$$

Ядра в SVM расширяют линейную модель классификации:

- ❶  $K(x, x') = (\langle x, x' \rangle + 1)^d$   
— полиномиальная разделяющая поверхность степени  $\leq d$ ;
- ❷  $K(x, x') = \sigma(\langle x, x' \rangle)$   
— нейронная сеть с заданной функцией активации  $\sigma(z)$   
( $K$  не при всех  $\sigma$  является ядром);
- ❸  $K(x, x') = \text{th}(k_1 \langle x, x' \rangle - k_0), \quad k_0, k_1 \geq 0$   
— нейросеть с сигмоидными функциями активации;
- ❹  $K(x, x') = \exp(-\gamma \|x - x'\|^2)$   
— сеть радиальных базисных функций (RBF ядро);

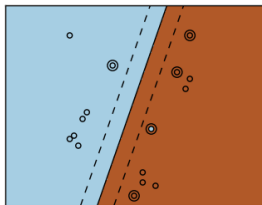
# Классификация с различными ядрами

Гиперплоскость в спрямляющем пространстве соответствует нелинейной разделяющей поверхности в исходном.

Примеры с различными ядрами  $K(x, x')$

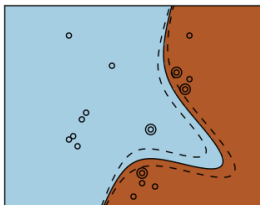
линейное

$$\langle x, x' \rangle$$



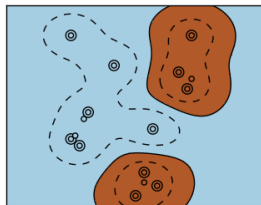
полиномиальное

$$(\langle x, x' \rangle + 1)^d, \quad d=3$$



гауссовское (RBF)

$$\exp(-\gamma \|x - x'\|^2)$$



---

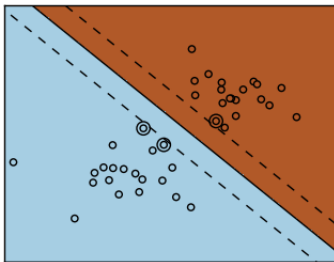
Пример из Python scikits learn: <http://scikit-learn.org/dev>

# Влияние константы $C$ на решение SVM

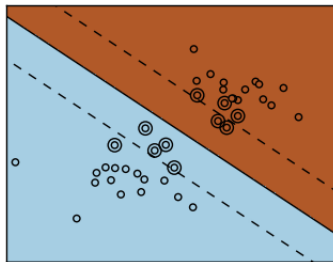
SVM — аппроксимация и регуляризация эмпирического риска:

$$\sum_{i=1}^{\ell} (1 - M_i(w, w_0))_+ + \frac{1}{2C} \|w\|^2 \rightarrow \min_{w, w_0}.$$

большой  $C$   
слабая регуляризация



малый  $C$   
сильная регуляризация



---

Пример из Python scikits learn: <http://scikit-learn.org/dev>

## Преимущества:

- Задача выпуклого квадратичного программирования имеет единственное решение.
- Выделяется множество опорных объектов.
- Имеются эффективные численные методы для SVM.
- Изящное обобщение на нелинейные классификаторы.

## Недостатки:

- Опорными объектами могут становиться выбросы.
- Нет отбора признаков в исходном пространстве  $X$ .
- Приходится подбирать константу  $C$ .