

# Лекция 7

## Метод опорных векторов

Габдуллин Р.А., Макаренко В.А.

МГУ им. М.В. Ломоносова

2 марта 2021

# Задача классификации

$X$  – множество объектов,

$Y$  – множество ответов:

- $|Y| = 2$  – двухклассовая (binary) классификация.
- $|Y| = K$  – множественная (multiclass) классификация.

$y : X \rightarrow Y$  – неизвестная зависимость.

**Дано:**

$\{x_1, x_2, \dots, x_\ell\} \subset X$  – обучающая выборка,

$y_i = y(x_i)$ ,  $i = 1, \dots, \ell$  – известные ответы.

**Найти:**

$a : X \rightarrow Y$  – решающая функция, приближающая  $y$  на всём  $X$ .

# Модель бинарной классификации

- Множество ответов:

$$Y = \{-1, 1\}.$$

- Семейство вещественных дискриминантных функций:

$$S = \{s(x, w) | w \in W\}.$$

- Семейство алгоритмов:

$$a(x, w) = \text{sign } s(x, w).$$

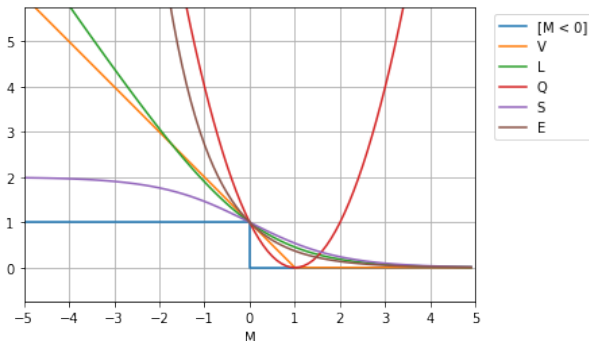
- Эмпирический риск:

$$Q(w, \mathbb{X}) = \sum_{i=1}^{\ell} [M(x_i, w) < 0] \equiv \sum_{i=1}^{\ell} [y_i \cdot s(x_i, w) < 0].$$

- Минимизация мажоранты эмпирического риска:

$$Q(w, \mathbb{X}) = \sum_{i=1}^{\ell} [M(x_i, w) < 0] \leq \sum_{i=1}^{\ell} \mathcal{L}(M(x_i, w)) \rightarrow \min_w.$$

# Мажоранты эмпирического риска



Часто используемые функции потерь  $\mathcal{L}$ :

- $V(M) = (1 - M)_+$
- $L(M) = \log_2(1 + e^{-M})$
- $Q(M) = (1 - M)^2$
- $S(M) = 2(1 + e^M)^{-1}$
- $E(M) = e^{-M}$

# Метод опорных векторов (support vector machine)

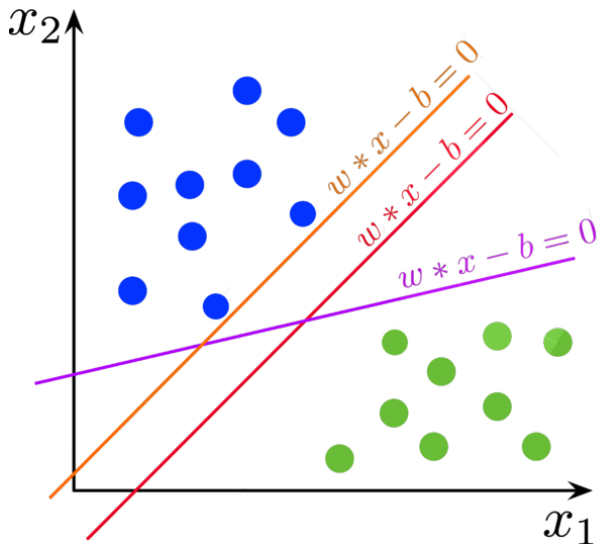


Рис.: Источник: [neerc.ifmo.ru](http://neerc.ifmo.ru)

# SVM: линейно разделимый случай

- Обучающая выборка:

$$X^\ell = \{(x_i, y_i)\}_{i=1}^\ell, \quad x_i \in \mathbb{R}^n, y_i \in \{-1, 1\}.$$

- Семейство алгоритмов:

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

- Отступ (margin) на  $i$ -м объекте:

$$M_i(w, b) = y_i(\langle w, x_i \rangle - b).$$

- Нормировка:

$$\min_{1 \leq i \leq \ell} M_i(w, b) = \min_{1 \leq i \leq \ell} |\langle w, x_i \rangle - b| = 1.$$

# SVM: линейно разделимый случай

- Цель – сделать расстояние от разделяющей гиперплоскости до ближайшего к ней объекта как можно больше:

$$\min_{1 \leq i \leq \ell} \frac{|\langle w, x_i \rangle - b|}{\|w\|} = \frac{1}{\|w\|} \rightarrow \max.$$

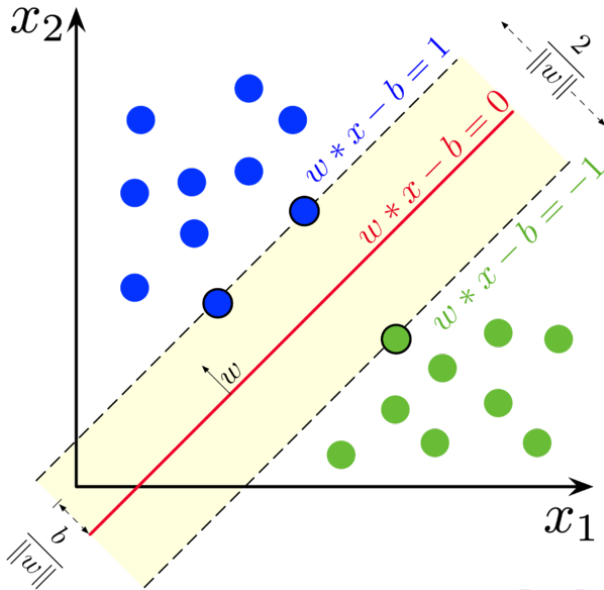
- Ширина разделяющей полосы (расстояние до ближайшего объекта положительного класса равно расстоянию до ближайшего объекта отрицательного класса):

$$\frac{2}{\|w\|}.$$

- Задача оптимизации:

$$\begin{cases} \min_{1 \leq i \leq \ell} M_i(w, b) = 1, \\ \frac{1}{2} \|w\|^2 \rightarrow \min. \end{cases} \iff \begin{cases} M_i(w, b) \geq 1, & 1 \leq i \leq \ell, \\ \frac{1}{2} \|w\|^2 \rightarrow \min. \end{cases}$$

# SVM: линейно разделимый случай





# SVM: линейно неразделимый случай

- Введем штрафы за попадание в разделяющую полосу или на территорию другого класса.
- Задача оптимизации:

$$\begin{cases} M_i(w, b) \geq 1 - \xi_i, & 1 \leq i \leq \ell, \\ \xi_i \geq 0, & 1 \leq i \leq \ell, \\ \frac{1}{2} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \rightarrow \min, \end{cases}$$

где  $C > 0$  – гиперпараметр.

- Эквивалентная задача безусловной оптимизации:

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

# Условия Каруша–Куна–Таккера

Задача нелинейного программирования:

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

Если  $x$  – точка локального минимума, то существуют такие множители  $\mu_i, \lambda_j$  ( $1 \leq i \leq m, 1 \leq j \leq k$ ), что для функции Лагранжа

$$L(x, \mu, \lambda) = f(x) + \sum_{i=1}^m \mu_i g_i(x) + \sum_{j=1}^k \lambda_j h_j(x)$$

выполняются условия

$$\begin{cases} \frac{\partial L}{\partial x} = 0, \\ g_i(x) \leq 0, h_j(x) = 0, & \text{(исходные ограничения)} \\ \mu_i \geq 0, & \text{(двойственные ограничения)} \\ \mu_i g_i(x) = 0 & \text{(условия дополняющей нежесткости)} \end{cases}$$

# Условия Каруша–Куна–Таккера в SVM

Функция Лагранжа:

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

Условия Каруша–Куна–Таккера:

$$\begin{cases} \frac{\partial L}{\partial w} = 0, & \frac{\partial L}{\partial b} = 0, & \frac{\partial L}{\partial \xi} = 0, \\ \xi_i \geq 0, & \lambda_i \geq 0, & \eta_i \geq 0, & M_i(w, b) \geq 1 - \xi_i, & 1 \leq i \leq \ell, \\ \lambda_i = 0 & \text{либо} & M_i(w, b) = 1 - \xi_i, & & 1 \leq i \leq \ell, \\ \eta_i = 0 & \text{либо} & \xi_i = 0, & & 1 \leq i \leq \ell. \end{cases}$$

# Условия Каруша–Куна–Таккера в SVM

Функция Лагранжа:

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

Продифференцируем и приравняем производные к нулю:

$$\frac{\partial L}{\partial w} = w - \sum_{i=1}^{\ell} \lambda_i y_i x_i = 0 \quad \Longleftrightarrow \quad w = \sum_{i=1}^{\ell} \lambda_i y_i x_i,$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^{\ell} \lambda_i y_i = 0 \quad \Longleftrightarrow \quad \sum_{i=1}^{\ell} \lambda_i y_i = 0,$$

$$\frac{\partial L}{\partial \xi_i} = -\lambda_i - \eta_i + C = 0 \quad \Longleftrightarrow \quad \lambda_i + \eta_i = C, \quad 1 \leq i \leq \ell.$$

# Условия Каруша–Куна–Таккера в SVM

Условия Каруша–Куна–Таккера:

$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i, & \sum_{i=1}^{\ell} \lambda_i y_i = 0, & \lambda_i + \eta_i = C, & 1 \leq i \leq \ell \\ \xi_i \geq 0, & \lambda_i \geq 0, & \eta_i \geq 0, & M_i(w, b) \geq 1 - \xi_i, & 1 \leq i \leq \ell, \\ \lambda_i = 0 & \text{либо} & M_i(w, b) = 1 - \xi_i, & 1 \leq i \leq \ell, \\ \eta_i = 0 & \text{либо} & \xi_i = 0, & 1 \leq i \leq \ell. \end{cases}$$

- Объект  $x_i$  называется опорным, если  $\lambda_i \neq 0$ .
- Можем разделить объекты на три типа:
  - 1  $\lambda_i = 0 \Rightarrow \eta_i = C, \xi_i = 0, M_i \geq 1$  – периферийные объекты,
  - 2  $0 < \lambda_i < C \Rightarrow 0 < \eta_i < C, \xi_i = 0, M_i = 1$  – опорные граничные объекты,
  - 3  $\lambda_i = C \Rightarrow \eta_i = 0, \xi_i > 0, M_i < 1$  – опорные объекты-нарушители.

# Двойственная задача

Подставляем в функцию Лагранжа полученные ограничения и приходим к двойственной задаче:

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

Решение прямой задачи выражается через решение двойственной:

$$\begin{cases} w = \sum_{i=1}^{\ell} \lambda_i y_i x_i, \\ b = \langle w, x_i \rangle - y_i. \end{cases}$$

Линейный классификатор принимает вид:

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

# Нелинейное обобщение, ядерный переход

Заменяем везде  $\langle x, x' \rangle$  нелинейной функцией  $K(x, x')$ , называемой ядром.

Примеры ядер:

- $K(x, x') = \langle x, x' \rangle^d$  – полиномиальное ядро.
- $K(x, x') = \exp(-\gamma \|x - x'\|^2)$  – сеть радиальных базисных функций.