



# Машинное обучение ФИБТ

## DS-поток

Лекция 7



# SVM

## Метод опорных векторов

Бинарная классификация



# SVM

Исходная задача:

*решать сложно*

$$\begin{cases} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i^+ \longrightarrow \min_{\theta, \theta_0, \xi} \\ Y_i(\theta^T X_i + \theta_0) \geq 1 - \xi_i \quad i = 1, \dots, n \end{cases}$$

Двойственная задача

*решать легко*

(квадр. прогр.)

$$\begin{cases} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j Y_i Y_j \langle X_i, X_j \rangle + \sum_{i=1}^n \lambda_i \longrightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C \quad i = 1, \dots, n; \\ \sum_{i=1}^n \lambda_i Y_i = 0 \end{cases}$$

Решение:  $\hat{\theta} = \sum_{i=1}^n \lambda_i Y_i X_i$ ;  $\hat{\theta}_0 = Y_i - \hat{\theta}^T X_i$  где  $i$  т.ч.  $\xi_i = 0$

$$\hat{y}(x) = \langle \hat{\theta}, x \rangle + \hat{\theta}_0$$



# SVM

Исходная задача

$$\begin{cases} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i^+ \longrightarrow \min_{\theta, \theta_0, \xi} \\ Y_i(\theta^T X_i + \theta_0) \geq 1 - \xi_i \quad i = 1, \dots, n \end{cases}$$

Преобразуем:

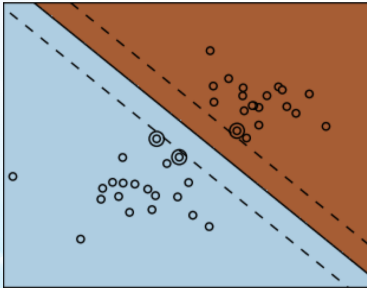
$$\sum_{i=1}^n (1 - M_i)^+ + \frac{1}{2C} \|\theta\|^2 \rightarrow \min_{\theta, \theta_0}$$

где  $M_i = Y_i(\theta^T X_i + \theta_0)$

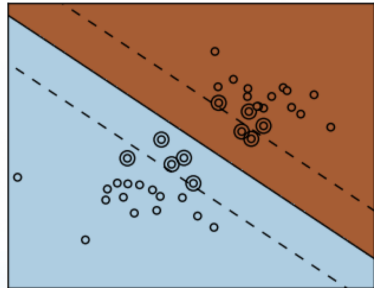
**Вывод:** SVM эквивалентен поиску линейного классификатора с ф-ией потерь  $L(y, z) = (1 - yz)^+ [hinge\ loss]$  и  $L_2$ -регуляризацией.

# SVM: влияние константы $C$

Большой  $C$   
Слабая регуляризация



Малый  $C$   
Сильная регуляризация



На практике небольшие изменения в значении  $C$   
не сильно влияют на вид классификатора.

⇒  $C$  можно подбирать по достаточно грубой сетке.



# SVM: ключевые моменты

- ▶ SVM — линейный классификатор с кусочно-линейной функцией потерь (hinge loss) и L2-регуляризацией.
- ▶ Придуман из соображений максим. зазора между классами.
- ▶ При линейно-разделимой выборке это означает максимизацию ширины разделяющей полосы.
- ▶ При линейно неразделимой выборке добавляется возможность попадания объектов в полосу и штрафы за эти попадания.
- ▶ Объекты в глубине классов не влияют на разделяющую гиперплоскость.



# Kernel trick

*Недостаток:*

Является линейной моделью  $\Rightarrow$  разделяющие полосы линейные.

Рассмотрим преобразование пространства объектов:  $\psi : \mathcal{X} \rightarrow \mathcal{H}$ ,  
где  $\mathcal{X} = \mathbb{R}^d$  — исходное пр-во,

$\mathcal{H}$  — гильбертово пр-во, возможно, бесконечномерное.

*Напоминание:* в гильбертовом пр-ве есть скалярное произведение.

Обозначим  $K(x_1, x_2) = \langle \psi(x_1), \psi(x_2) \rangle$

*Наблюдения:*

1.  $\theta \in \mathcal{X} \Rightarrow$  к нему тоже применимо преобразование  $\psi$ .
2. используем только скалярные произведения.



# SVM

Исходная задача:

*решать сложно*

$$\begin{cases} \frac{1}{2} \langle \theta, \theta \rangle + C \sum_{i=1}^n \xi_i^+ \longrightarrow \min_{\theta, \theta_0, \xi} \\ Y_i (\langle \theta, X_i \rangle + \theta_0) \geq 1 - \xi_i \quad i = 1, \dots, n \end{cases}$$

Двойственная задача

*решать легко*

(квадр. прогр.)

$$\begin{cases} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j Y_i Y_j \langle X_i, X_j \rangle + \sum_{i=1}^n \lambda_i \longrightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C \quad i = 1, \dots, n; \\ \sum_{i=1}^n \lambda_i Y_i = 0 \end{cases}$$

Решение:  $\hat{\theta} = \sum_{i=1}^n \lambda_i Y_i X_i$ ;  $\hat{\theta}_0 = Y_i - \langle \hat{\theta}, X_i \rangle$  где  $i$  т.ч.  $\xi_i = 0$

$$\hat{y}(x) = \langle \hat{\theta}, x \rangle + \hat{\theta}_0$$





# SVM

Исходная задача:

*решать сложно*

$$\begin{cases} \frac{1}{2} \langle \theta, \theta \rangle + C \sum_{i=1}^n \xi_i^+ \longrightarrow \min_{\theta, \theta_0, \xi} \\ Y_i (\langle \theta, X_i \rangle + \theta_0) \geq 1 - \xi_i \quad i = 1, \dots, n \end{cases}$$

Двойственная задача

*решать легко*

(кв. прогр.)

$$\begin{cases} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j Y_i Y_j \langle X_i, X_j \rangle + \sum_{i=1}^n \lambda_i \longrightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C \quad i = 1, \dots, n; \\ \sum_{i=1}^n \lambda_i Y_i = 0 \end{cases}$$

Решение:  $\hat{\theta} = \sum_{i=1}^n \lambda_i X_i$ ,  $\hat{\theta}_0 = \frac{1}{n} \sum_{i=1}^n Y_i \langle \hat{\theta}, X_i \rangle + 1$

$$\hat{y}(x) = \sum_{i=1}^n \lambda_i Y_i \langle X_i, x \rangle + \left( Y_\ell - \sum_{i=1}^n \lambda_i Y_i \langle X_i, X_\ell \rangle \right), \text{ т.ч. } X_\ell \text{ на границе.}$$



# SVM

Исходная задача:

*решать сложно*

$$\begin{cases} \frac{1}{2} K(\theta, \theta) + C \sum_{i=1}^n \xi_i^+ \rightarrow \min_{\theta, \theta_0, \xi} \\ Y_i(K(\theta, X_i) + \theta_0) \geq 1 - \xi_i \quad i = 1, \dots, n \end{cases}$$

Двойственная задача

*решать легко*

(квадр. прогр.)

$$\begin{cases} -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j Y_i Y_j K(X_i, X_j) + \sum_{i=1}^n \lambda_i \rightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C \quad i = 1, \dots, n; \\ \sum_{i=1}^n \lambda_i Y_i = 0 \end{cases}$$

Решение:  $\hat{\theta} = \frac{1}{n} \sum_{i=1}^n \lambda_i Y_i X_i$ ,  $\hat{\theta}_0 = \frac{1}{n} \sum_{i=1}^n \lambda_i Y_i$

$$\hat{y}(x) = \sum_{i=1}^n \lambda_i Y_i K(X_i, x) + \left( Y_\ell - \sum_{i=1}^n \lambda_i Y_i K(X_i, X_\ell) \right), \text{ т.ч. } X_\ell \text{ на границе}$$



# Kernel trick

*Наблюдение:* не нужно знать саму  $\psi$ .

**Определение:**

$K(x_1, x_2)$  — ядро, если  $\exists \psi : \mathcal{X} \rightarrow \mathcal{H}$  т.ч.  $K(x_1, x_2) = \langle \psi(x_1), \psi(x_2) \rangle$ ,  
где  $\mathcal{H}$  — некоторое гильбертово пространство.

*Замечание.* Не путать с ядром из KDE.

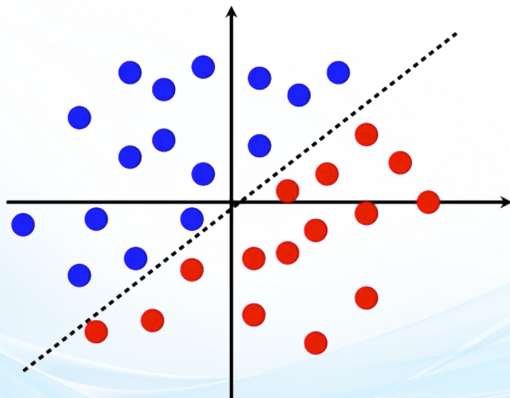
**Теорема:**

Функция  $K(x_1, x_2)$  является ядром  $\iff$

- ▶  $K(x_1, x_2)$  симметрична:  $K(x_1, x_2) = K(x_2, x_1)$ ;
- ▶  $K(x_1, x_2)$  неотрицательно определена:  
т.е. для любых  $x_1, \dots, x_n \in \mathcal{X}$  матрица  $(K(x_i, x_j))_{ij}$  неотр. определена.

# Kernel trick

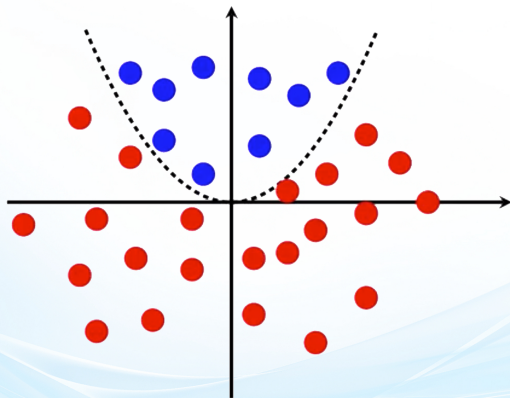
Линейное ядро



$$K(x_1, x_2) = \langle x_1, x_2 \rangle$$

# Kernel trick

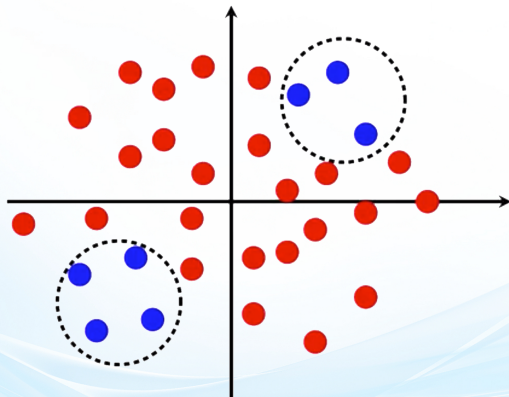
Полиномиальное ядро



$$K(x_1, x_2) = (\gamma \langle x_1, x_2 \rangle + r)^d$$

# Kernel trick

Радиальное ядро



$$K(x_1, x_2) = e^{-\gamma \langle x_1 - x_2 \rangle^2}$$



## Kernel trick

*Почему это ядра?*

Рассмотрим пример квадратичного ядра, т.е.  $K(x, z) = \langle x, z \rangle^2$ .

Пусть  $\mathcal{X} = \mathbb{R}^2$ , т.е.  $x = (x_1, x_2)$ ,  $z = (z_1, z_2)$

Разложим  $K(x, z)$ :

$$\begin{aligned} K(x, z) &= \langle x, z \rangle^2 = \langle (x_1, x_2), (z_1, z_2) \rangle^2 = \\ &= (x_1 z_1 + x_2 z_2)^2 = x_1^2 z_1^2 + x_2^2 z_2^2 + 2x_1 z_1 x_2 z_2 = \\ &= \left\langle \begin{pmatrix} x_1^2, x_2^2, \sqrt{2}x_1 x_2 \end{pmatrix}, \begin{pmatrix} z_1^2, z_2^2, \sqrt{2}z_1 z_2 \end{pmatrix} \right\rangle \end{aligned}$$

Таким образом,  $\mathcal{H} = \mathbb{R}^3$ ,  $\psi(x) = (x_1^2, x_2^2, \sqrt{2}x_1 x_2)$

Линейная поверхность в  $\mathcal{H}$  соответствует квадратичной в  $\mathcal{X}$ .

Пространство  $\mathcal{H}$  называется **спрямляющим**.



# Недостатки SVM

- ▶ Нужно подбирать  $C$ .
- ▶ Нет общих рекомендаций для выбора  $K(x, z)$ .
- ▶ Работает только для бинарной классификации.





# SVM

## Метод опорных векторов

Многоклассовый случай



# Многоклассовый SVM

Теперь  $Y_i \in \{1, \dots, K\}$ .

Пусть в выборке имеется константный признак  
т.е. не нужно явно указывать сдвиг  $\theta_0$ .

**Идея:** отделяем класс  $k$  от всех остальных по знач.  $\langle \theta_k, x \rangle$ .  
где  $\theta_k$  — вектор параметров.

Итоговое предсказание имеет вид

$$\hat{y}(x) = \arg \max_{k \in \{1, \dots, K\}} \langle \theta_k, x \rangle.$$

Двуклассовый случай	$\begin{cases} \frac{1}{2} \ \theta_1\ ^2 \rightarrow \min_{\theta_1} \\ Y_i \langle \theta_1, X_i \rangle \geq 1, \quad i = 1, \dots, n \end{cases}$
Многоклассовый	$\begin{cases} \frac{1}{2} \sum_{k=1}^K \ \theta_k\ ^2 \rightarrow \min_{\theta} \\ \langle \theta_{Y_i}, X_i \rangle - \langle \theta_k, X_i \rangle \geq 1, \quad i = 1..n; \quad k \in \{1..K\} \setminus \{Y_i\} \end{cases}$



# Многоклассовый SVM

## Общий случай

Двуклассовый:

$$\begin{cases} \frac{1}{2} \|\theta\|^2 + C \sum_{i=1}^n \xi_i^+ \longrightarrow \min_{\theta, \theta_0, \xi} \\ Y_i(\theta^T X_i + \theta_0) \geq 1 - \xi_i \quad i = 1, \dots, n \end{cases}$$

Многоклассовый:

$$\begin{cases} \frac{1}{2} \sum_{k=1}^K \|\theta_k\|^2 + C \sum_{i=1}^n \xi_i^+ \longrightarrow \min_{\theta, \xi} \\ \langle \theta_{Y_i}, X_i \rangle - \langle \theta_k, X_i \rangle \geq 1 - \xi_i, \quad i = 1..n; \quad k \in \{1..K\} \setminus \{Y_i\} \end{cases}$$

# Эквивалентная функция потерь

Рассмотрим следующую **функцию потерь**:

$$\mathcal{L}(X_i) = \max_k \left\{ \langle \theta_k, X_i \rangle + 1 - I\{k = Y_i\} \right\} - \langle \theta_{Y_i}, X_i \rangle$$

Выражение, по которому берется max:

- ▶ Если  $k = Y_i$ , то оно равно  $\langle \theta_k, X_i \rangle$ , иначе оно равно  $\langle \theta_k, X_i \rangle + 1$

Получаем

- ▶ Если оценка за верный класс больше оценок за остальные классы хотя бы на единицу, т.е.  $\forall k : \langle \theta_{Y_i}, X_i \rangle > \langle \theta_k, X_i \rangle + 1$ ,  
 $\Rightarrow$  Максимум достигается на  $k = Y_i$   
 $\Rightarrow \mathcal{L}(X_i) = \langle \theta_{Y_i}, X_i \rangle - \langle \theta_{Y_i}, X_i \rangle = 0$
- ▶ Иначе, т.е.  $\exists k : \langle \theta_{Y_i}, X_i \rangle < \langle \theta_k, X_i \rangle + 1$   
 $\Rightarrow$  Максимум достигается на  $k \neq Y_i$   
 $\Rightarrow \mathcal{L}(X_i) = \langle \theta_k, X_i \rangle + 1 - \langle \theta_{Y_i}, X_i \rangle > 0$

**Вывод:** Штрафуем как за неверный ответ на объекте, так и за попадание в разделяющую полосу



# SVM

## Метод опорных векторов

История



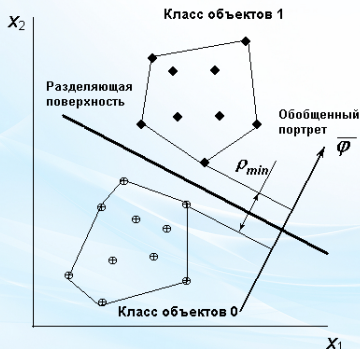
# Метод обобщенного портрета

Предложен в 1960-е Вапником и Червоненкисом.

**Первая версия:** Отделение отделеение гиперплоскостью точек одного класса на сфере от всего остального.

**Вторая версия:** Разделение гиперплоскостью двух классов на сфере.

**Третья версия:** Разделение гиперплоскостью двух классов в пр-ве.

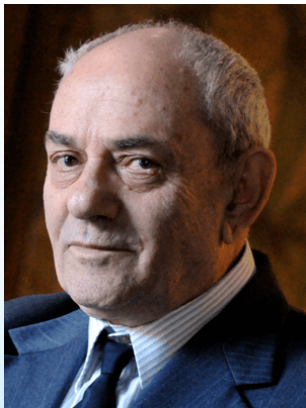


**SVM:** предложен Вапником в 1990-е

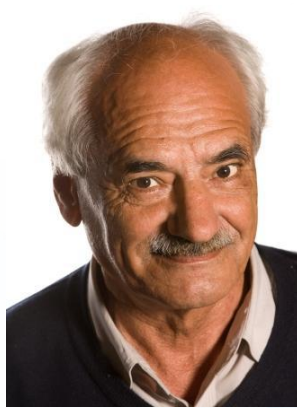
Ключевые отличия:

1. Ошибки  $\xi_i$
2. Ядра.

В середине 2000-х был популярным.



Вапник Владимир Наумович  
род. 1936  
учился в Узбекистане



Червоненкис Алексей Яковлевич  
1938-2014  
выпускник МФТИ, 1961