

# Введение в анализ данных

Лекция 15

Метод опорных векторов

Евгений Соколов

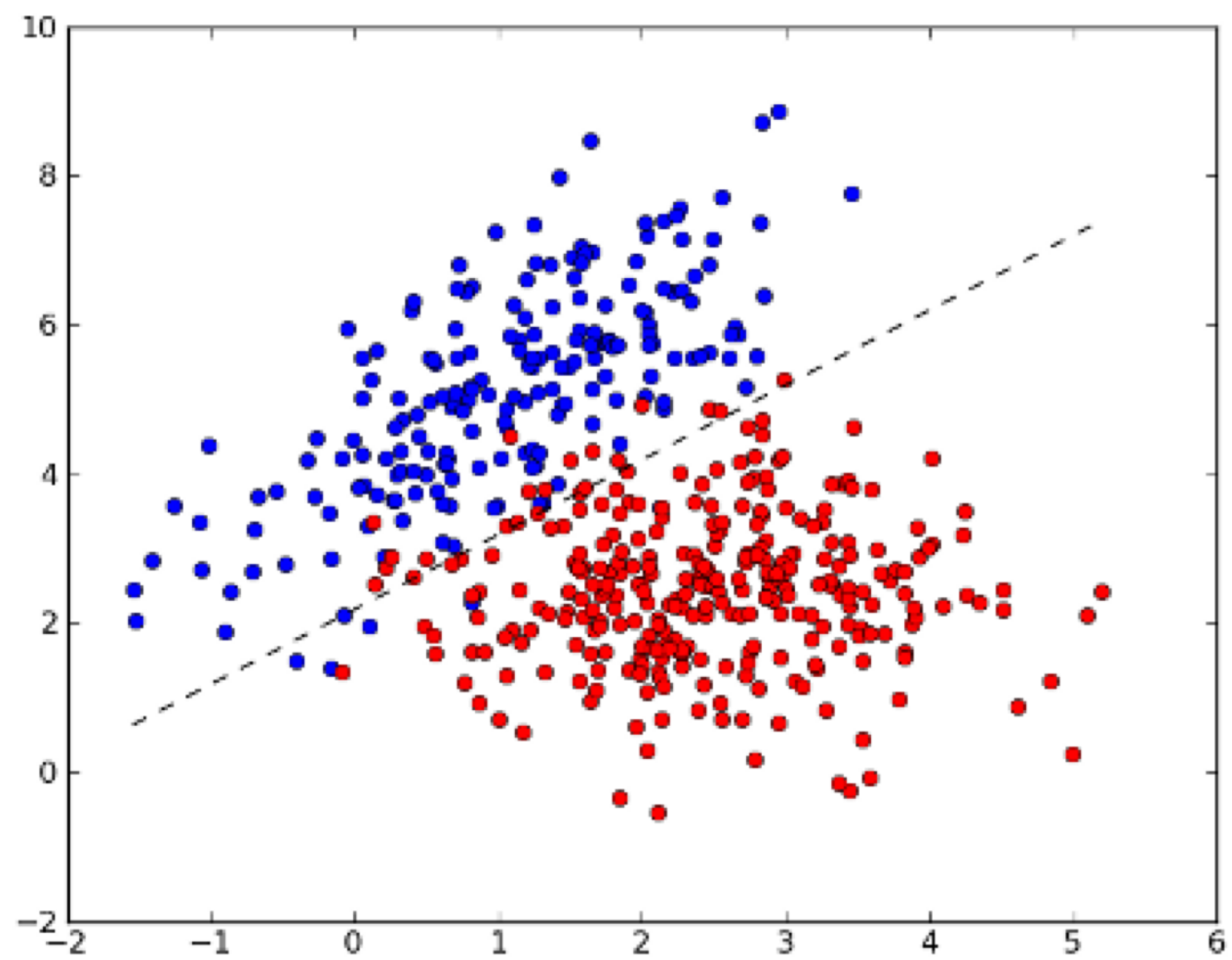
[esokolov@hse.ru](mailto:esokolov@hse.ru)

НИУ ВШЭ, 2019

# Логистическая регрессия

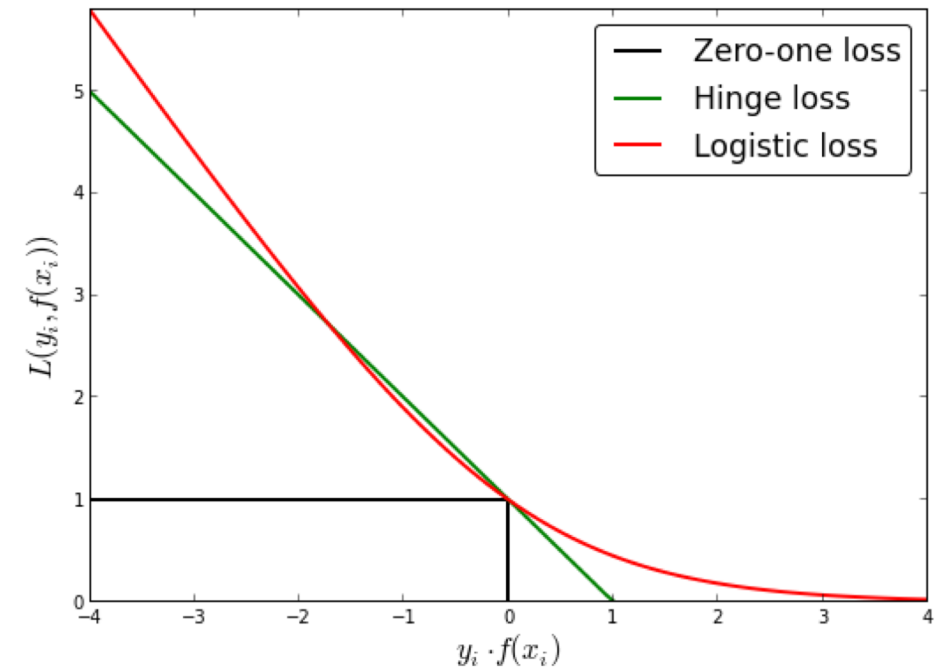
- Линейная модель классификации:  $a(x) = \text{sign } \langle w, x \rangle$
- Позволяет оценивать вероятности:  $\pi(x) = \sigma(\langle w, x \rangle)$

$$\frac{1}{\ell} \sum_{i=1}^{\ell} \log_2(1 + \exp(-y_i \langle w, x_i \rangle)) \rightarrow \min_w$$



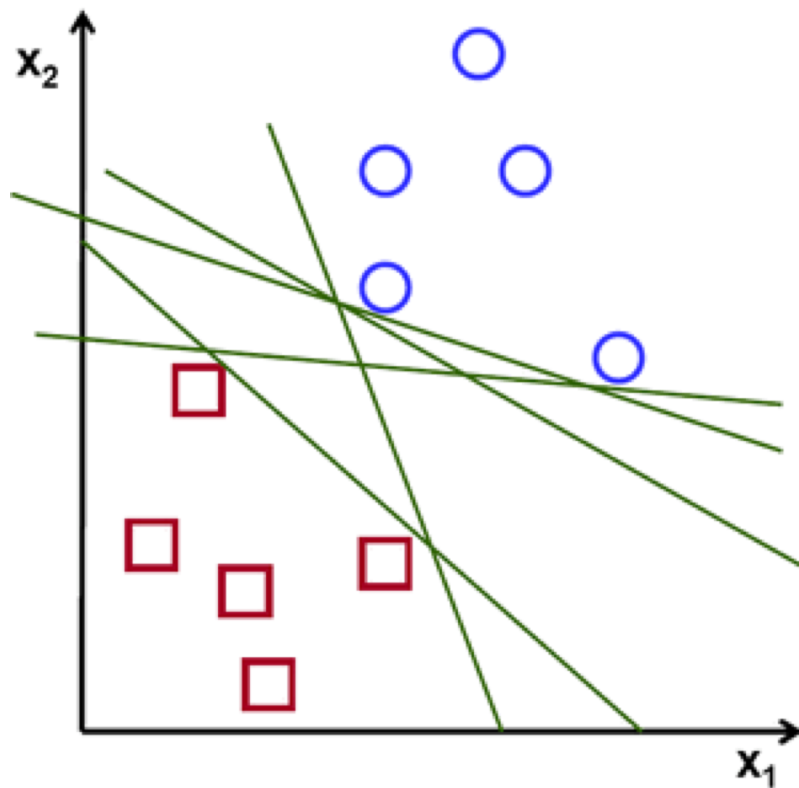
# Смысл штрафа

- Отступ должен быть как можно больше
- Логично, поскольку вероятность должна быть как можно ближе к нулю или единице
- А если забыть про вероятности и требовать только корректной классификации?



# Идея максимизации отступа

Объекты должны быть как можно дальше от разделяющей прямой



# Идея максимизации отступа

Предположение:

выборку можно идеально разделить линейным классификатором

Т.е. существует такое  $w$ , что для всех объектов обучающей выборки

$$y_i(\langle w, x_i \rangle + w_0) > 0$$

# Идея максимизации отступа

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

Подготовка: если поделить  $w$  и  $w_0$  на положительное число, то ответы классификатора не поменяются

Пример:

$$\text{sign}(10 * x_1 + 4 * x_2 + 2) = \text{sign}(5 * x_1 + 2 * x_2 + 1)$$

# Идея максимизации отступа

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

Подготовка: если поделить  $w$  и  $w_0$  на положительное число, то ответы классификатора не поменяются

Поделим так, что выполнено **условие нормировки**

$$\min_{x \in X} |\langle w, x \rangle + w_0| = 1$$

(на обучающей выборке минимальный модуль прогноза равен 1)



# Геометрия линейного классификатора

- Расстояние от точки до гиперплоскости  $\langle w, x \rangle + w_0 = 0$ :

$$\frac{|\langle w, x \rangle + w_0|}{\|w\|}$$

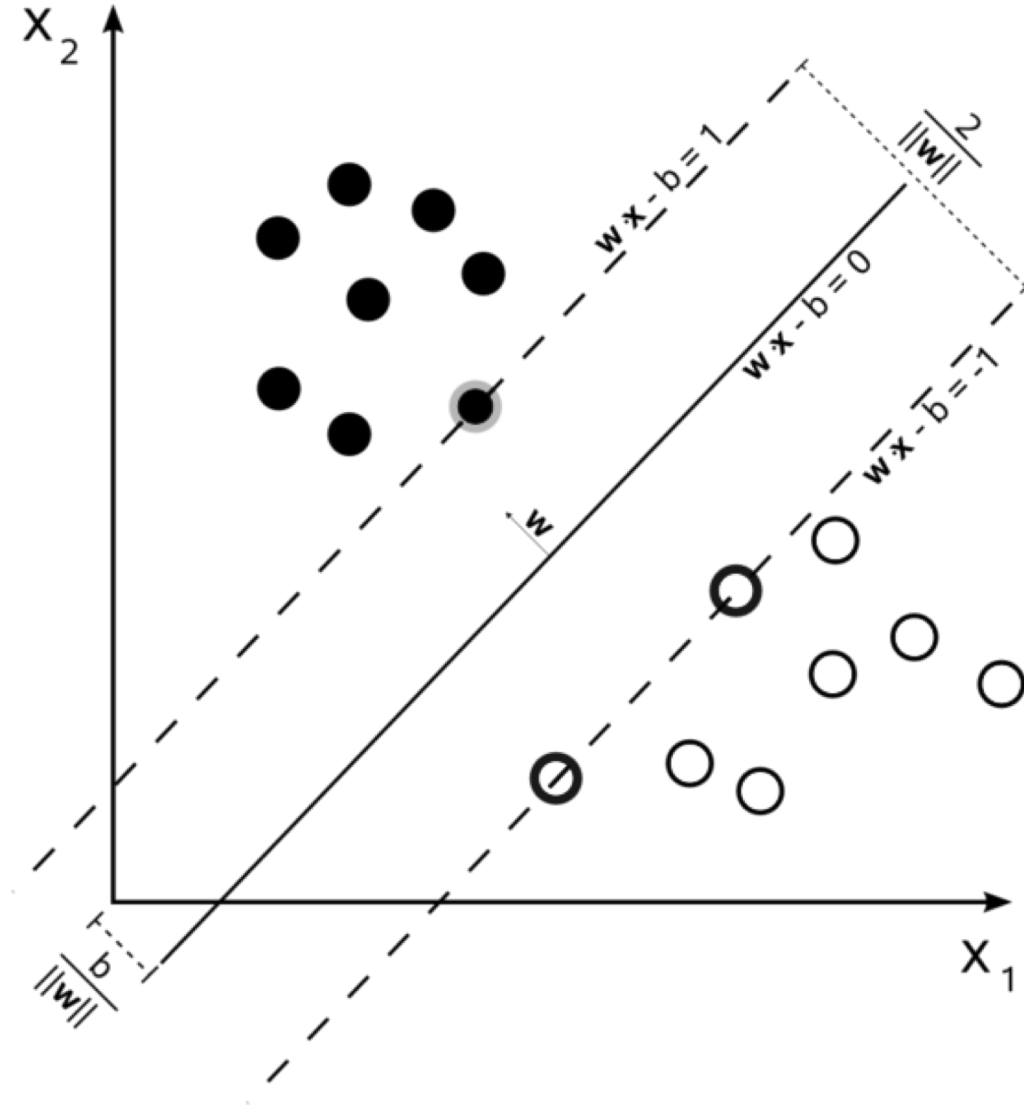
- Чем больше  $\langle w, x \rangle$ , тем дальше объект от разделяющей гиперплоскости

# Идея максимизации отступа

Минимальное расстояние от объекта обучающей выборки до разделяющей гиперплоскости:

$$\min_{x \in X} \frac{|\langle w, x \rangle + w_0|}{\|w\|} = \frac{1}{\|w\|} \min_{x \in X} |\langle w, x \rangle + w_0| = \frac{1}{\|w\|}$$

То есть ширина разделяющей полосы равна  $\frac{2}{\|w\|}$



# Задача максимизации ширины полосы

$$\begin{cases} \|w\|^2 \rightarrow \min_{w, w_0} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 \end{cases}$$

Неравенство равносильно двум другим неравенствам:

- $y_i(\langle w, x_i \rangle + w_0) > 0$  — все объекты должны правильно классифицироваться
- $|\langle w, x \rangle + w_0| \geq 1$  — условие нормировки

# Задача максимизации ширины полосы

$$\begin{cases} \|w\|^2 \rightarrow \min_{w, w_0} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 \end{cases}$$

Для решения существуют специальные методы оптимизации

# Задача максимизации ширины полосы

- Всё это время мы изучали задачу, которая никогда не встречается!
- Вряд ли настоящие данные будут линейно разделимыми

# Задача максимизации ширины полосы

$$\begin{cases} \|w\|^2 \rightarrow \min_{w, w_0} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 \end{cases}$$

- Условие невозможно выполнить для всех объектов на линейно неразделимой выборке
- Сделаем условие более мягким

# Задача максимизации ширины полосы

$$\begin{cases} \|w\|^2 \rightarrow \min_{w, w_0} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases}$$

- Теперь мы разрешаем отступу на некоторых объектах быть меньше единицы
- Что не так с этой задачей?



# Задача максимизации ширины полосы

$$\begin{cases} \|w\|^2 \rightarrow \min_{w, w_0} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases}$$

- Теперь мы разрешаем отступу на некоторых объектах быть меньше единицы
- Что не так с этой задачей?
- Можно взять  $\xi_i = +\infty$ , и тогда подойдёт решение  $w = 0$

# Метод опорных векторов (SVM)

$$\begin{cases} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \rightarrow \min_{w, w_0, \xi_i} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases}$$

Параметр  $C$  отвечает за баланс между шириной полосы и качеством классификации:

- При больших  $C$  стараемся не допускать ошибок
- При малых  $C$  разрешаем не обращать внимание на много объектов

# Метод опорных векторов (SVM)

$$\begin{cases} \|w\|^2 + C \sum_{i=1}^{\ell} \xi_i \rightarrow \min_{w, w_0, \xi_i} \\ y_i(\langle w, x_i \rangle + w_0) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{cases}$$

Последние два условия равносильны одному равенству:

$$\xi_i = \max(0, 1 - y_i(\langle w, x_i \rangle + w_0))$$

# Метод опорных векторов (SVM)

Последние два условия равносильны одному равенству:

$$\xi_i = \max(0, 1 - y_i(\langle w, x_i \rangle + w_0))$$

Эквивалентная задача:

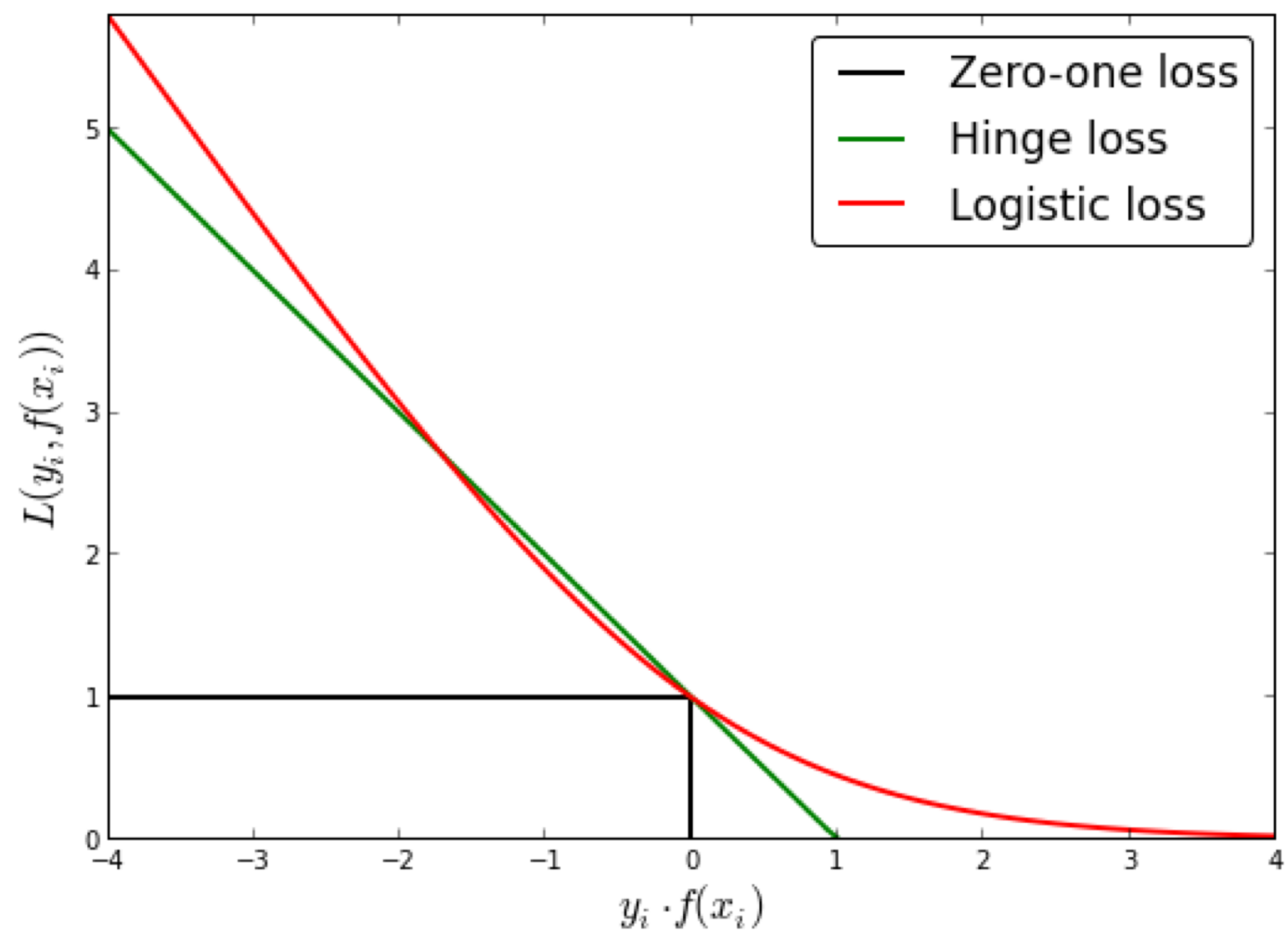
$$\|w\|^2 + C \sum_{i=1}^{\ell} \max(0, 1 - y_i(\langle w, x_i \rangle + w_0)) \rightarrow \min_{w, w_0}$$

# Метод опорных векторов (SVM)

$$\|w\|^2 + C \sum_{i=1}^{\ell} \max(0, 1 - y_i(\langle w, x_i \rangle + w_0)) \rightarrow \min_{w, w_0}$$

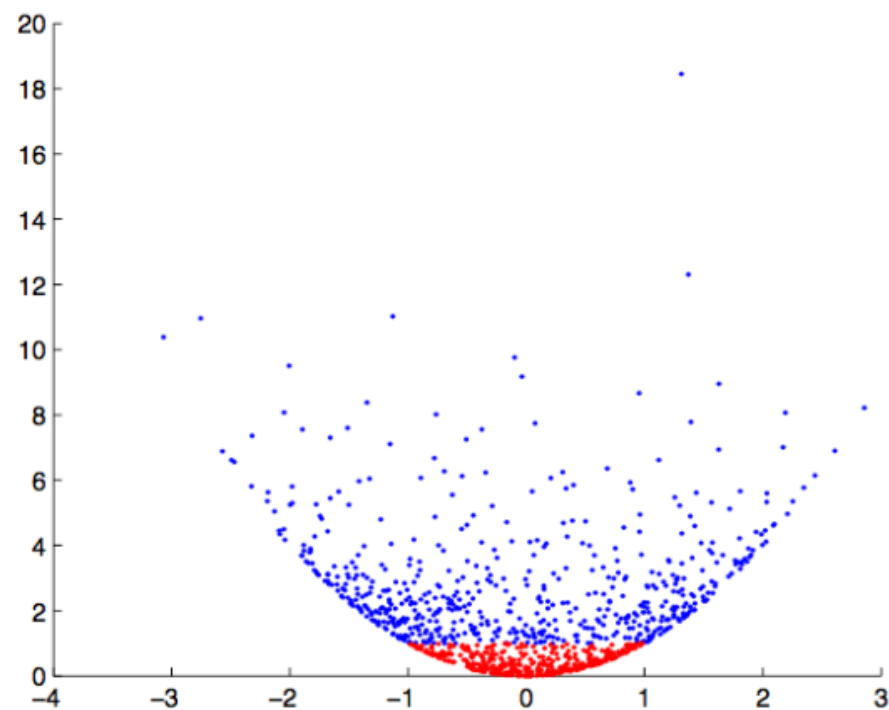
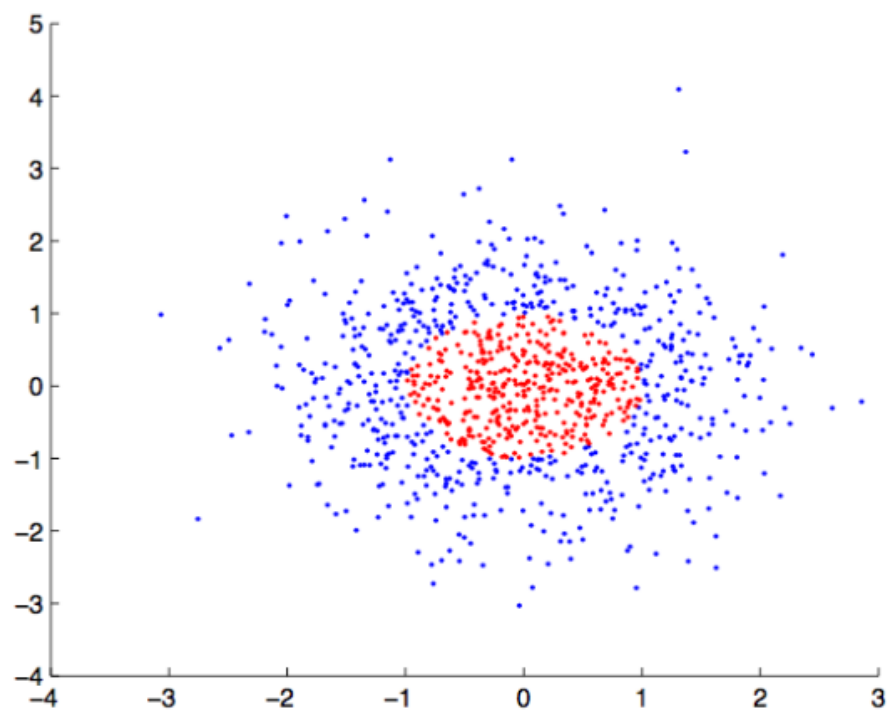
Регуляризатор

Ошибка с функцией потерь  
 $L(M) = \max(0, 1 - M)$



Как сделать линейную модель  
нелинейной?

# Нелинейные признаки





# Нелинейные признаки

- Можно добавлять к исходным признакам их нелинейные преобразования
- Пример:  $x_i^2$ ,  $x_i x_j$
- Если исходных признаков 100, то парных — несколько тысяч
- Если исходных признаков 1000, то парных — сотни тысяч
- Проблемы с памятью и производительностью

# Нелинейные признаки

- Допустим, мы полностью переходим к парным признакам:

$$\phi(x) = (x_i x_j)_{i,j=1}^{\ell}$$

- Скалярное произведение:

$$\langle \phi(x), \phi(z) \rangle = \sum_{i=1}^{\ell} \sum_{j=1}^{\ell} x_i x_j z_i z_j = \sum_{i=1}^{\ell} x_i z_i \sum_{j=1}^{\ell} x_j z_j = \langle x, z \rangle^2$$

- При некоторых нелинейных признаках скалярные произведения считаются легко и без дополнительной памяти

Ядровой переход

# Двойственная задача

$$\left\{ \begin{array}{l} \sum_{i=1}^{\ell} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell, \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{array} \right.$$

Что общего у этой задачи и SVM?

# Двойственная задача

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

Что общего у этой задачи и SVM?

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i$$

# Двойственная задача

$$w = \sum_{i=1}^{\ell} \lambda_i y_i x_i$$

Модель:

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

# Двойственная задача

$$\left\{ \begin{array}{l} \sum_{i=1}^{\ell} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \lambda_i \lambda_j y_i y_j \langle x_i, x_j \rangle \rightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell, \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{array} \right.$$

- Важно: задача зависит от объектов только через их скалярные произведения!
- Можно заменить на скалярные произведения, соответствующие нелинейным признакам!

# Двойственная задача

$$\left\{ \begin{array}{l} \sum_{i=1}^{\ell} \lambda_i - \frac{1}{2} \sum_{i,j=1}^{\ell} \lambda_i \lambda_j y_i y_j K(x_i, x_j) \rightarrow \max_{\lambda} \\ 0 \leq \lambda_i \leq C, \quad i = 1, \dots, \ell, \\ \sum_{i=1}^{\ell} \lambda_i y_i = 0. \end{array} \right.$$

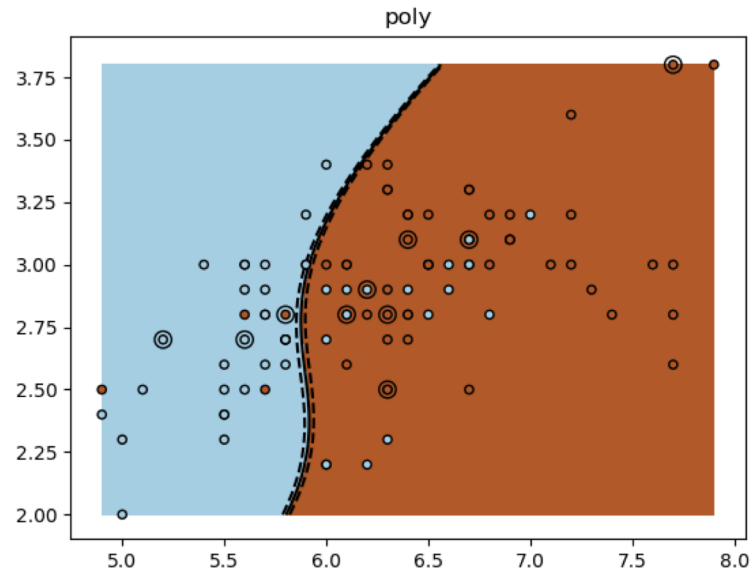
Здесь  $K(x, z)$  — **ядро**, новое скалярное произведение



# Примеры ядер

$$K(x, z) = (\langle x, z \rangle + R)^m$$

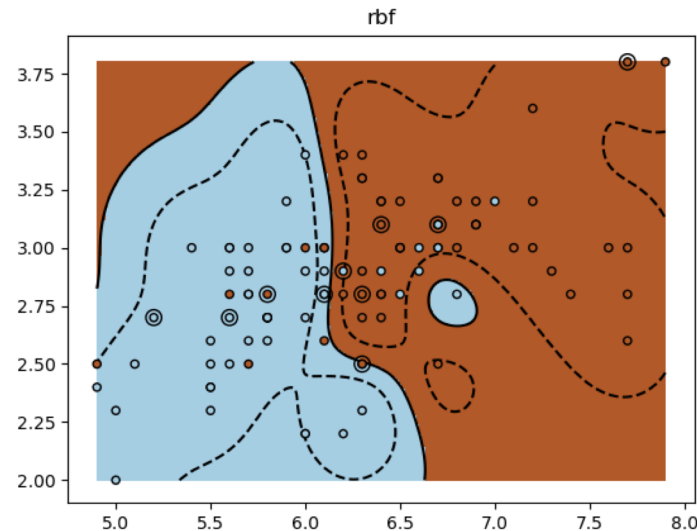
- Полиномиальное ядро
- Соответствует добавлению всех мономов степени не выше  $m$



# Примеры ядер

$$K(x, z) = \exp\left(-\frac{\|x-z\|^2}{2\sigma^2}\right)$$

- Гауссово ядро
- Соответствует добавлению **всех** мономов



# Метод опорных векторов

- Требуется правильной классификации объектов и наличия небольшого отступа
- Может быть преобразован так, что использует только скалярные произведения
- Ядра позволяют без дополнительных затрат подменить исходные признаковые описания на нелинейные