



Тверской
государственный
технический
университет

Интеллектуальные информационные системы

Метод опорных векторов

Материалы курса доступны по ссылке:

<https://github.com/AndreyShpigar/ML-course>

2024 г.

Метод опорных векторов в задачах классификации

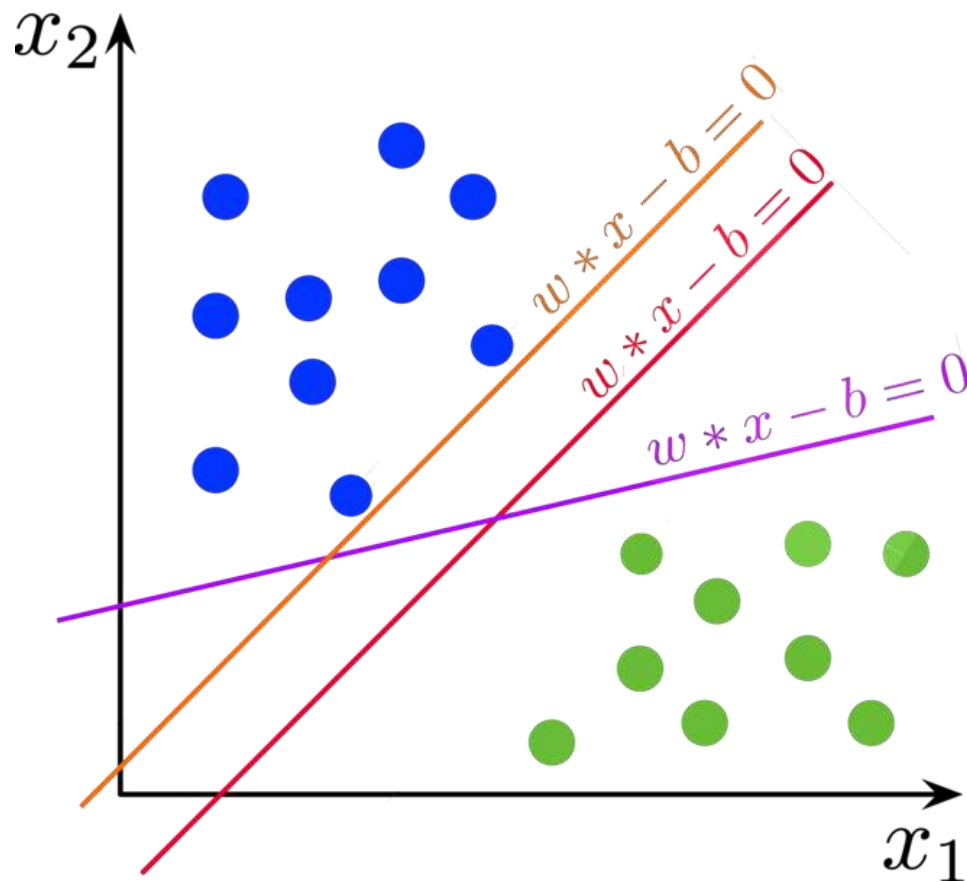
Задача бинарной классификации: $X = \mathbb{R}^n$,
 $Y = \{-1, +1\}$

Линейный классификатор:

$$a(x) = \operatorname{sign} \left(\sum_{j=1}^n w_j x^j - w_0 \right) = \operatorname{sign}(\langle w, x \rangle - w_0)$$

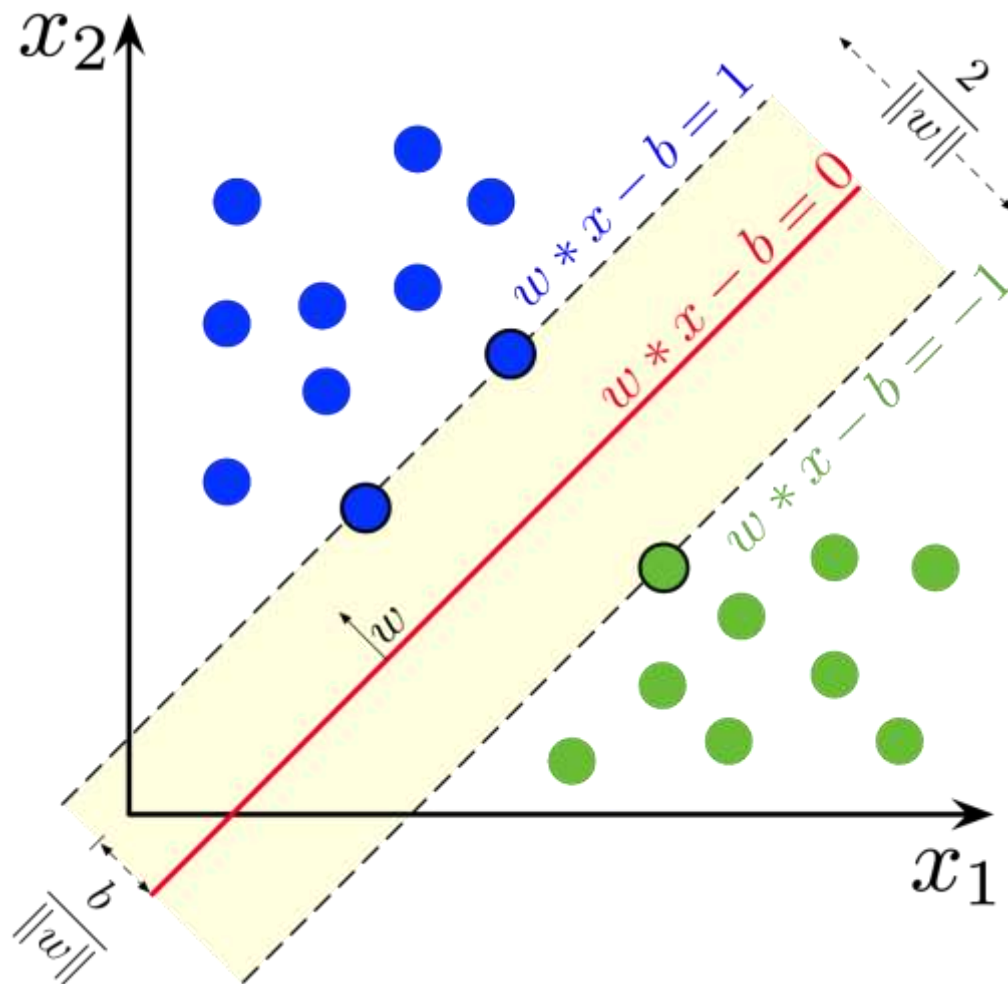
Предположим, что выборка линейно разделима:

$$Q(w, w_0) = \sum_{i=1}^l [y_i(\langle w, x_i \rangle - w_0) < 0] = 0$$



Ширина полосы:

$$\left\langle (x_+ - x_-), \frac{w}{\|w\|} \right\rangle = \frac{\langle w, x_+ \rangle - \langle w, x_- \rangle}{\|w\|} = \frac{(w_0 + 1) - (w_0 - 1)}{\|w\|} = \frac{2}{\|w\|}$$



Для линейно разделимой выборки:

Классификация с жестким зазором (*hard margin*)

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

Задача квадратичного программирования - минимизировать квадратичный функционал при линейных ограничениях

Обобщение для линейно неразделимой выборки:

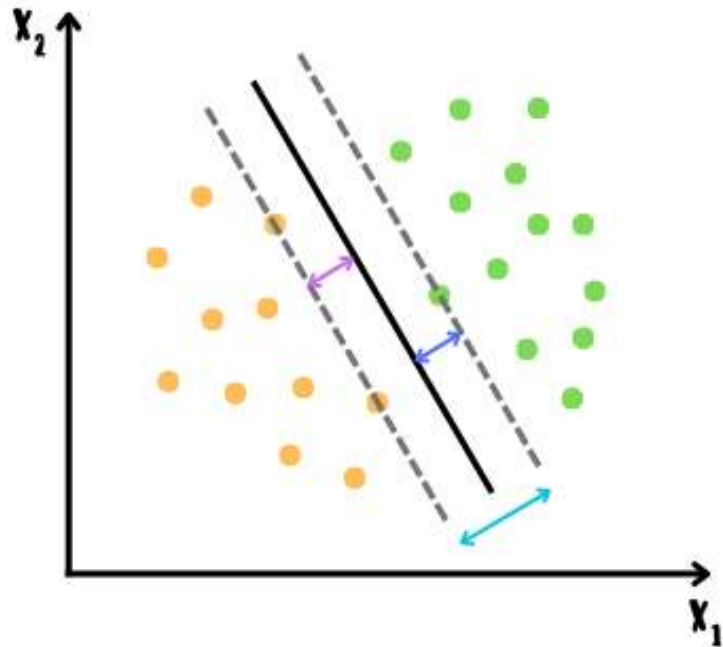
Классификация с мягким зазором (*soft margin*)

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

ξ_i — величина ошибки на объектах x_i

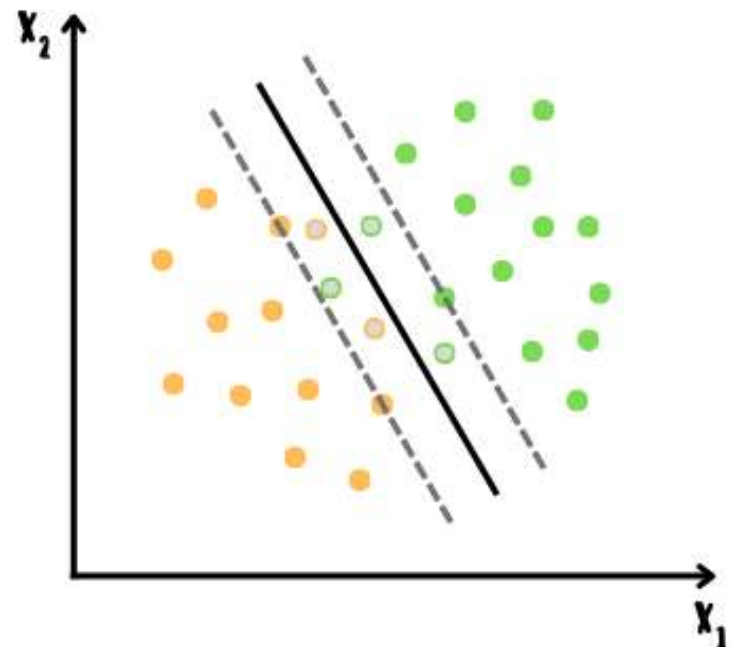
C — штраф за суммарную ошибку

hard margin



- margin
- margin
- total margin
- support vectors

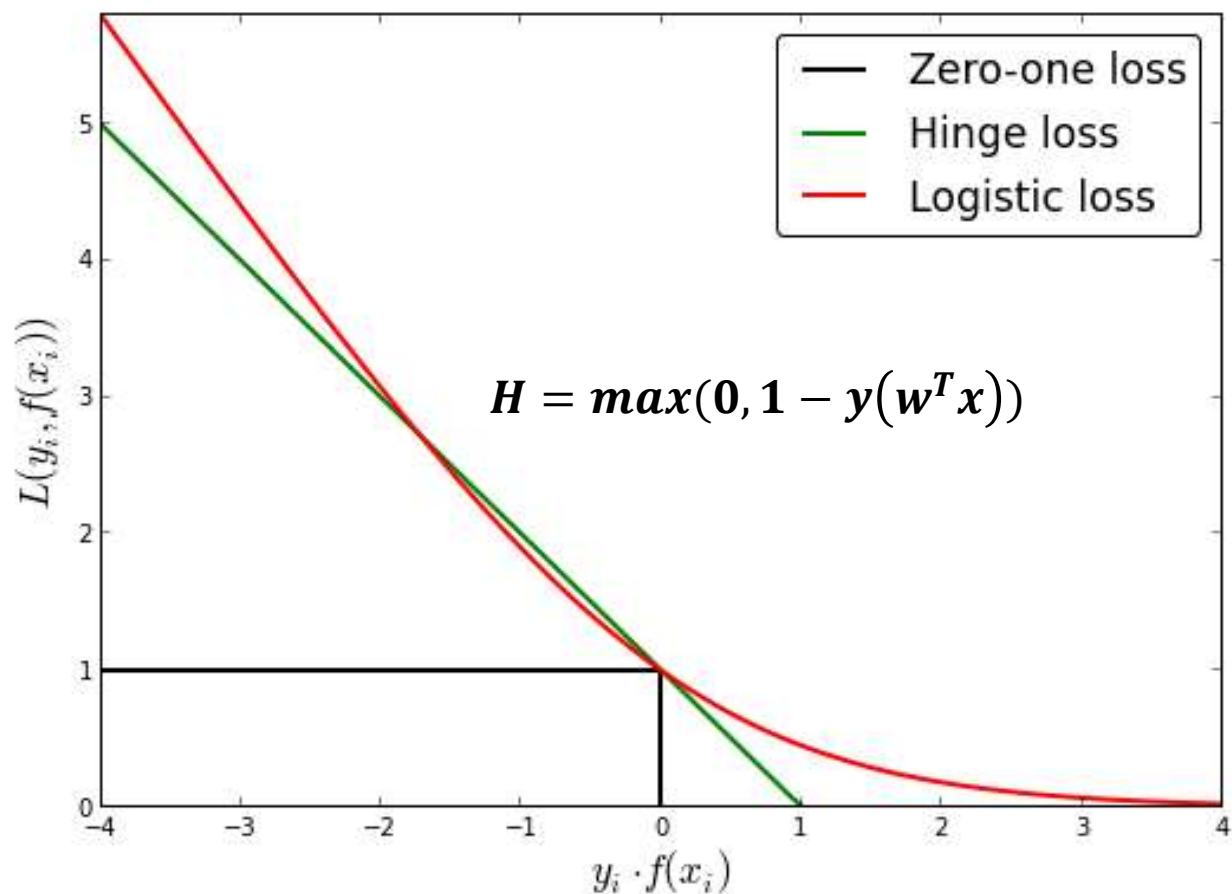
soft margin



- samples outside the support vectors
- samples outside the support vectors and the hyperplane

Эквивалентная задача безусловной оптимизации:

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



Инструмент для решения задачи оптимизации с ограничениями равенства и неравенства –
условия Каруша-Куна-Таккера

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

По теореме ККТ эта задача эквивалентна двойственной задаче поиска седловой точки функции Лагранжа

Если x – точка локального минимума, то существуют множители $\mu_i, i = 1, \dots, m$ и $\lambda_j, j = 1, \dots, k$:

$$\left\{ \begin{array}{l} \frac{\partial \mathcal{L}}{\partial x} = \mathbf{0}, \quad \mathcal{L}(x; \mu, \lambda) = f(x) + \sum_{i=1}^m \mu_i g_i(x) + \sum_{j=1}^k \lambda_j h_j(x) \\ g_i(x) \leq \mathbf{0}; \quad h_j(x) = \mathbf{0} \\ \mu_i \geq \mathbf{0} \\ \mu_i \vee g_i(x) = \mathbf{0} \end{array} \right.$$

Функция Лагранжа: $\mathcal{L}(w, w_0, \xi; \lambda, \eta) =$

$$\frac{1}{2} \|w\|^2 - \sum_{i=1}^l \lambda_i (M_i(w, w_0) - 1) - \sum_{i=1}^l \xi_i (\lambda_i + \eta_i - C)$$

Продифференцируем функцию Лагранжа и приравняем нулю производные:

$$\frac{\partial \mathcal{L}}{\partial \mathbf{w}} = \mathbf{0} \Rightarrow \mathbf{w} = \sum_{i=1}^l \lambda_i y_i \mathbf{x}_i$$

$$\frac{\partial \mathcal{L}}{\partial w_0} = 0 \Rightarrow \sum_{i=1}^l \lambda_i y_i = 0$$

$$\frac{\partial \mathcal{L}}{\partial \xi_i} = 0 \Rightarrow \eta_i + \lambda_i = C, \quad i = 1, \dots, l$$

Вектор весов \mathbf{w} выражается через λ_i :

$$\mathbf{w} = \sum_{i=1}^l \lambda_i y_i \mathbf{x}_i$$

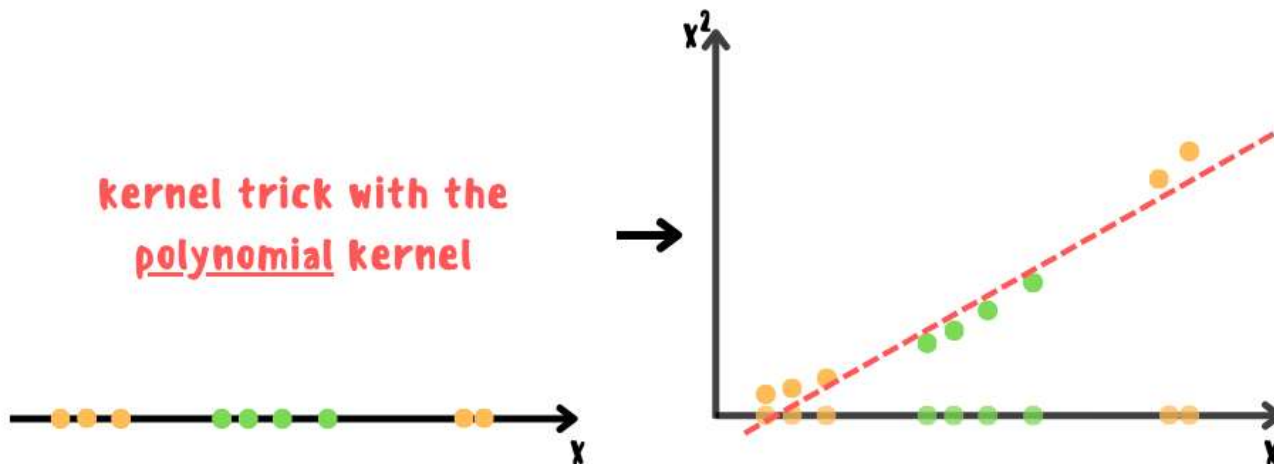
Если $\lambda_i = 0$, то решение задачи никак не зависит от объекта обучающей выборки \mathbf{x}_i (эти объекты можно назвать неинформативными)

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

Нелинейное обобщение SVM

Идея: переход от исходного пространства признаков описаний объектов X к новому пространству H (спрямляющему пространству) с помощью некоторого преобразования $\psi: X \rightarrow H$

Kernel trick (ядерный трюк) – замена скалярного произведения векторов n -й степени заменяется на их произведение в степени n : $\psi(a)^T \psi(b) = (a^T b)^n$



- Линейное ядро:

$$K(x, x^T) = \langle x, x^T \rangle$$

- Квадратичное ядро:

$$K(x, x^T) = \langle x, x^T \rangle^2$$

- Полиномиальное ядро:

$$K(x, x^T) = (\langle x, x^T \rangle + 1)^d$$

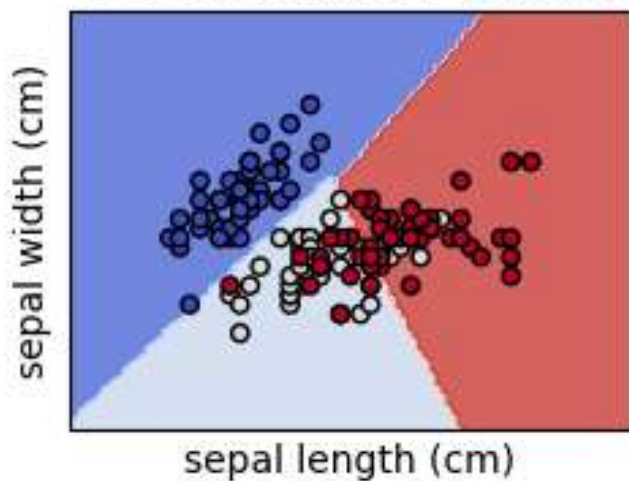
- Гауссовское RBF ядро (радиально-базисная функция):

$$K(x, x^T) = \exp(-\gamma \|x - x^T\|^2)$$

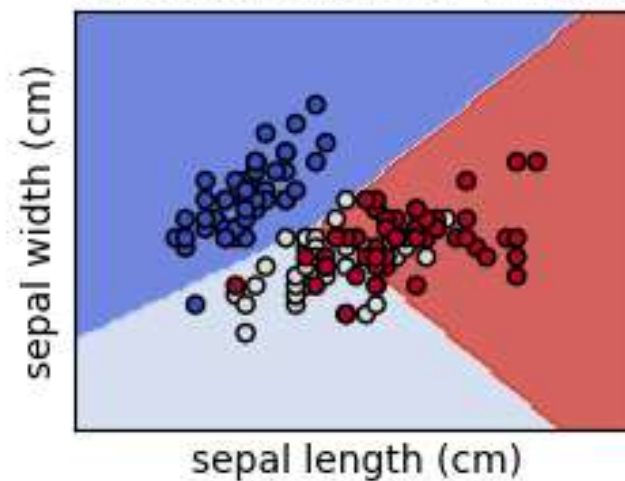
- Сигмоидальное ядро:

$$K(x, x^T) = \tanh(k_1 \langle x, x^T \rangle - k_0), \quad k_0, k_1 \geq 0$$

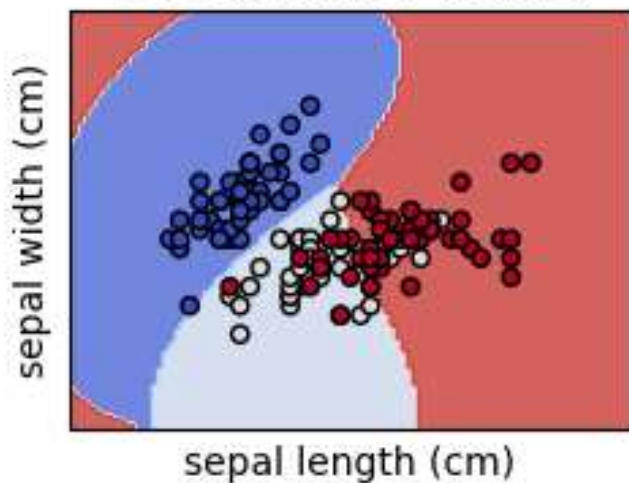
SVC with linear kernel



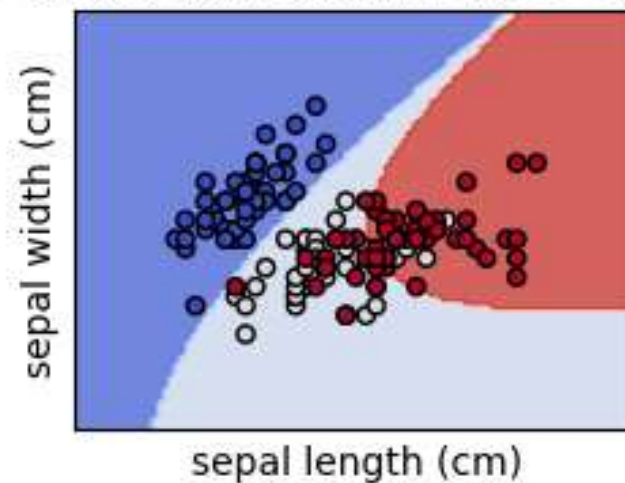
LinearSVC (linear kernel)

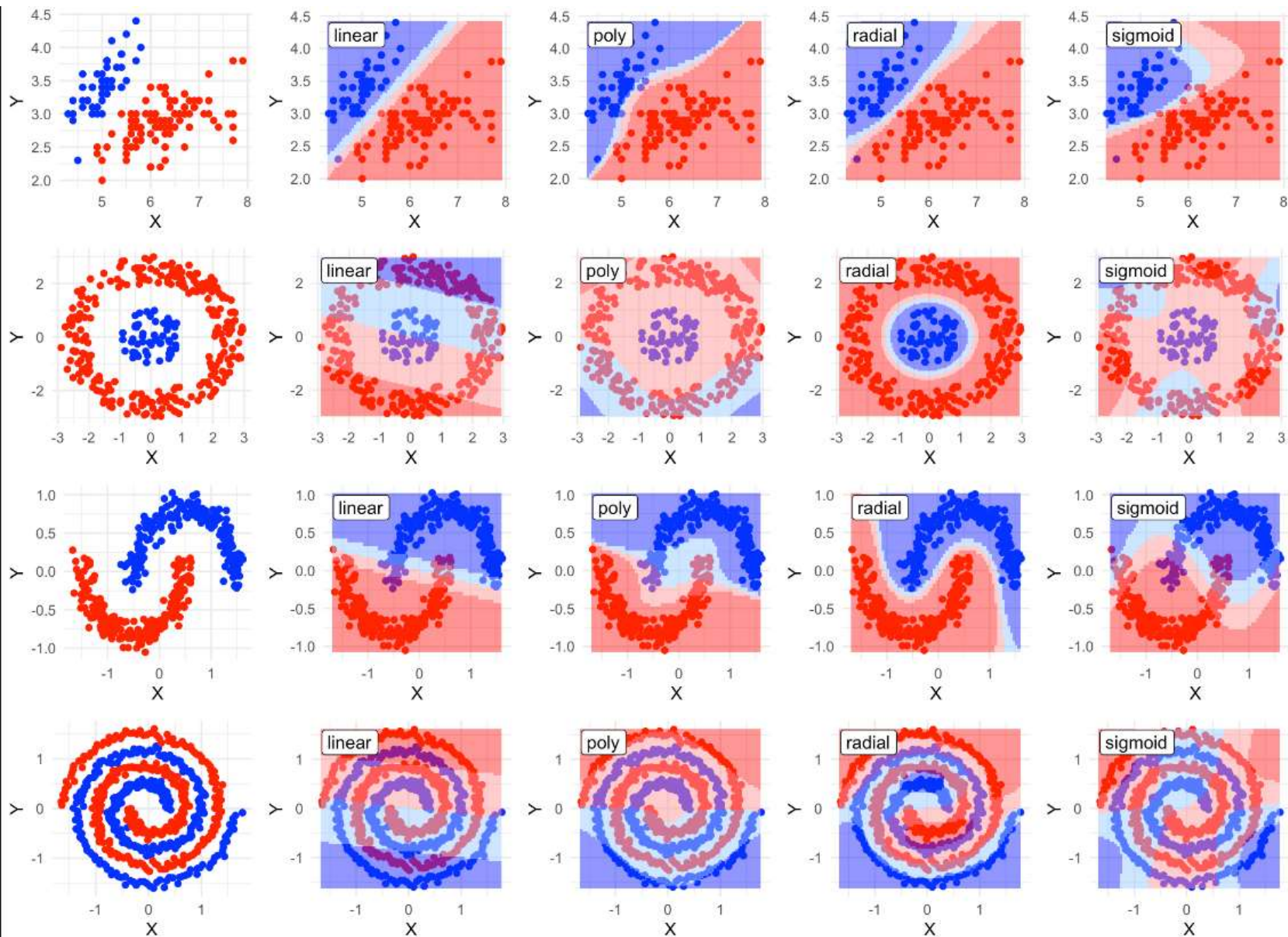


SVC with RBF kernel



SVC with polynomial (degree 3) kernel





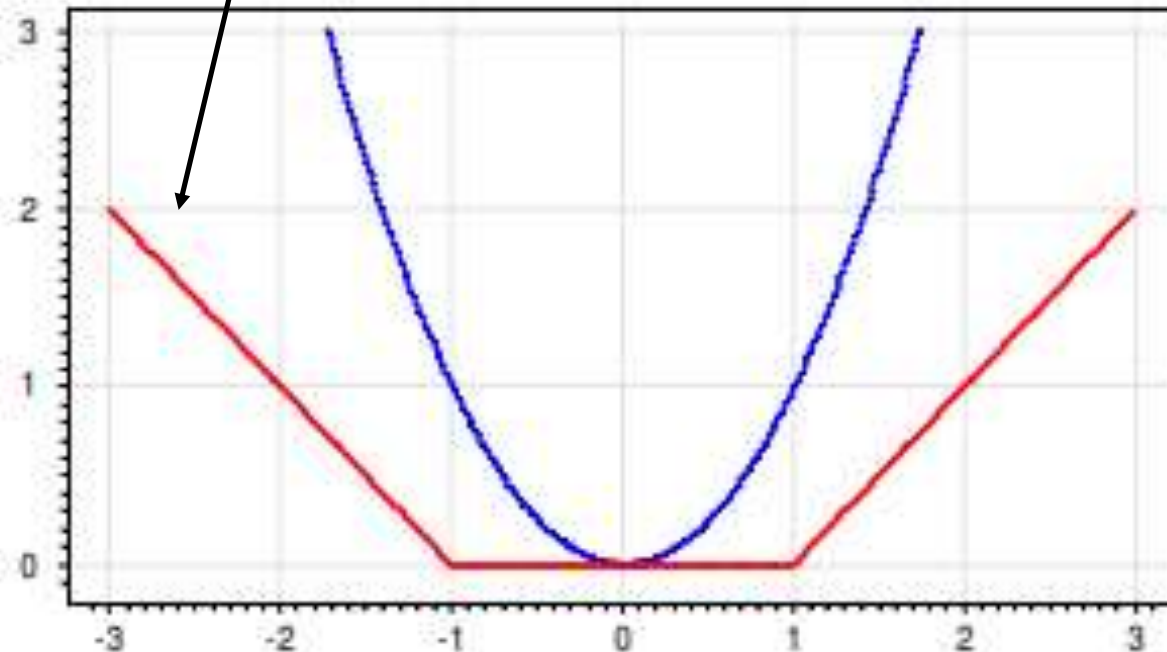
Метод опорных векторов в задачах регрессии

Модель регрессии:

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

Функция потерь: $L(\varepsilon) = (|\varepsilon| - \delta)_+$

Кусочно-линейная функция ε -чувствительности



Смысл функции потерь – если мы находимся на расстоянии не более чем δ от правильного ответа, то потери нет, за такой ответ не штрафует. Если дальше, то штраф увеличивается линейно.

Постановка задачи (с L2-регуляризацией):

$$\sum_{i=1}^l (|\langle \mathbf{w}, \mathbf{x}_i \rangle - w_0 - y_i| - \delta)_+ + \frac{1}{2C} \|\mathbf{w}\|^2 \rightarrow \min_{\mathbf{w}, w_0}$$

Достоинства SVM:

- ✓ Задача выпуклого квадратичного программирования имеет единственное решение – методы оптимизации существенно более эффективны
- ✓ Принцип оптимальной разделяющей гиперплоскости приводит к максимизации разделяющей полосы между классами – более уверенная классификация
- ✓ Возможность обработки многомерных данных без предварительного преобразования

Недостатки SVM:

- × Неустойчивость к шуму, выбросы в обучающей выборке могут стать опорными объектами-нарушителями и непосредственно влиять на построение разделяющей гиперплоскости
- × Не существует общих методов построения ядер и спрямляющих пространств
- × Необходимо подбирать параметр C в случае линейной неразделимости