

## Дополнительное задание 2

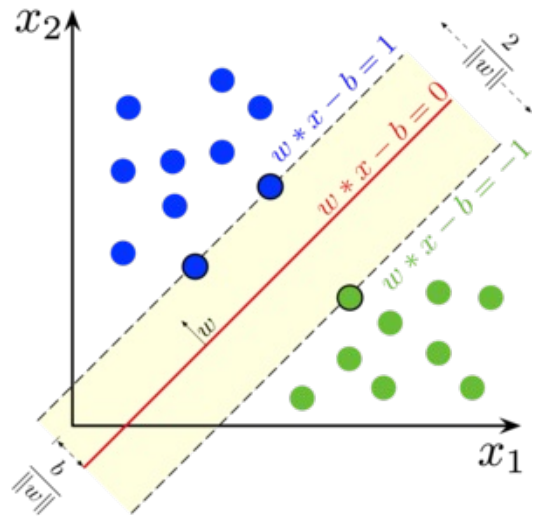
Рассмотрим метод опорных векторов (SVM). Это алгоритм машинного обучения, применяемый для задач линейной и нелинейной классификации и регрессии. Идея метода состоит в создании гиперплоскости, которая наилучшим образом разделяет объекты выборки. Алгоритм строится из следующего предположения: что чем больше расстояние (margin) между разделяющей гиперплоскостью и объектами различных классов, тем менее возможна ошибка классификации.

Рассмотрим бинарную задачу классификации:  $X \rightarrow Y$ , где  $X = \mathbb{R}^n$ ,  $Y = \{-1, 1\}$ . Пусть задана обучающая выборка пар:  $S = (\vec{x}_i, y_i)$ . Необходимо построить алгоритм классификации  $a(\vec{x}): X \rightarrow Y$ .

Любая гиперплоскость может быть записана как множество точек, удовлетворяющих уравнению  $w^T x - b = 0$ , где  $w$  - вектор нормали к гиперплоскости.

Пусть у нас есть гиперплоскость, которая делит данные на 2 класса  $C_1$  и  $C_2$ , тогда эти объекты располагаются по разные стороны от гиперплоскости, то есть

$$\begin{cases} \langle \vec{w}, \vec{x} \rangle - b > 0, & \forall x \in C_1 \\ \langle \vec{w}, \vec{x} \rangle - b < 0, & \forall x \in C_2 \end{cases}$$



Отступ (зазор, margin) — свойство, показывающее насколько объект "погружён" в свой класс, насколько типичным представителем класса он является.

Чем меньше отступ, тем ближе объекты находятся к гиперплоскости и тем самым вероятность ошибки тоже больше. Если отступ отрицательный, то классификатор допустил ошибку. Для линейного классификатора отступ определяется уравнением:

$$M(\vec{w}) = y_i(\langle \vec{w}, \vec{x}_i \rangle - b)$$

### Линейный SVM:

Предположим что выборка линейно разделима, тогда в качестве алгоритма классификации будем использовать линейный пороговый классификатор  $a: X \rightarrow Y$  вида

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

Так как линейная выборка разделима, то существует гиперплоскость, у которой отступ отступ до каждого объекта положителен и тогда построим такую гиперплоскость, что объекты обучающей выборки находились на наибольшем расстоянии от неё. Эти крайние гиперплоскости можно описать уравнениями:  $w^T x - b = 1$  и  $w^T x - b = -1$ .

Для того, чтобы определить ширину полосы, введём  $x_+$  и  $x_-$  (произвольных опорных объекта классов 1, -1). Тогда ширина полосы выражается как проекция вектора  $x_+$ ,  $x_-$  на нормаль к гиперплоскости  $w$ .

$$\begin{aligned} \frac{\langle \vec{x}_+ - \vec{x}_-, \vec{w} \rangle}{\|\vec{w}\|} &= \frac{\langle \vec{x}_+, \vec{w} \rangle - \langle \vec{x}_-, \vec{w} \rangle - b + b}{\|\vec{w}\|} = \frac{(+1)(\langle \vec{x}_+, \vec{w} \rangle - b) + (-1)(\langle \vec{x}_-, \vec{w} \rangle - b)}{\|\vec{w}\|} = \\ &= \frac{M_+(\vec{w}, b) + M_-(\vec{w}, b)}{\|\vec{w}\|} = \frac{2}{\|\vec{w}\|} \rightarrow \max \Rightarrow \|\vec{w}\| \rightarrow \min \end{aligned}$$

Задача в терминах квадратичного программирования:

$$(w, w) \rightarrow \min$$

$$M_i(x_i) \geq 1$$

### Нелинейный SVM:

В жизни линейный SVM не встречается особо, поэтому немного изменим немного правила игры и позволим объектам попадать на область другого класса и введем некий штраф за ошибку. Тогда можно записать задачу квадратичного программирования таким образом:

$$0.5 * (w, w) + C * \sum(\epsilon_i) \rightarrow \min$$

$$M_i(x_i) \geq 1 - \epsilon_i$$

$$\epsilon_i \geq 0$$

В качестве примера нелинейного SVM можно рассмотреть следующую окружность:

Как видно провести линию не получится, тогда воспользуемся нелинейными функциями для преобразования данных в более высокоразмерное пространство. В данном случае  $2D \rightarrow 3D$ , и тогда уже можем поделить.

