

Машинное обучение

Лекция 5

Линейные модели (Продолжение)

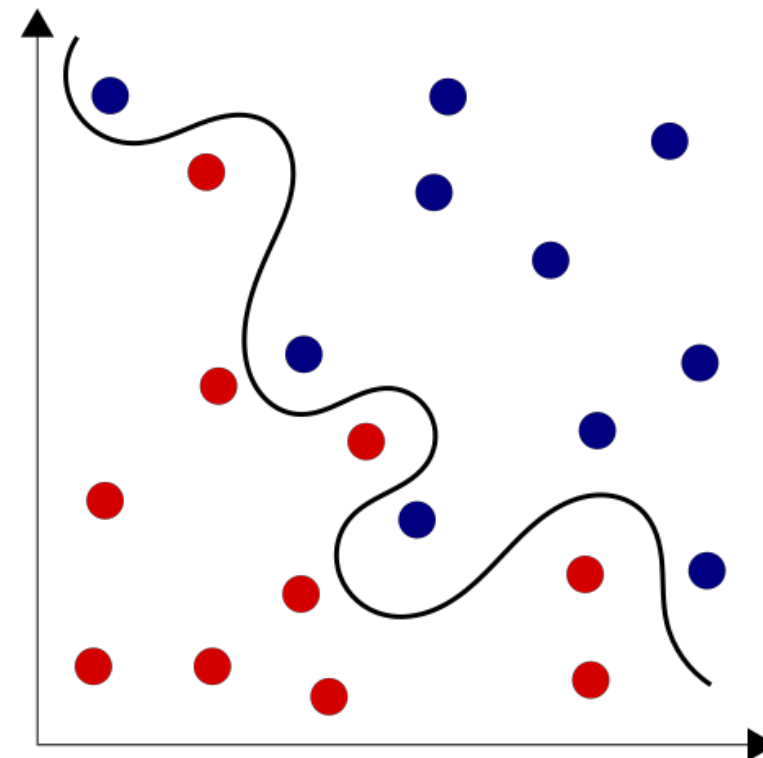
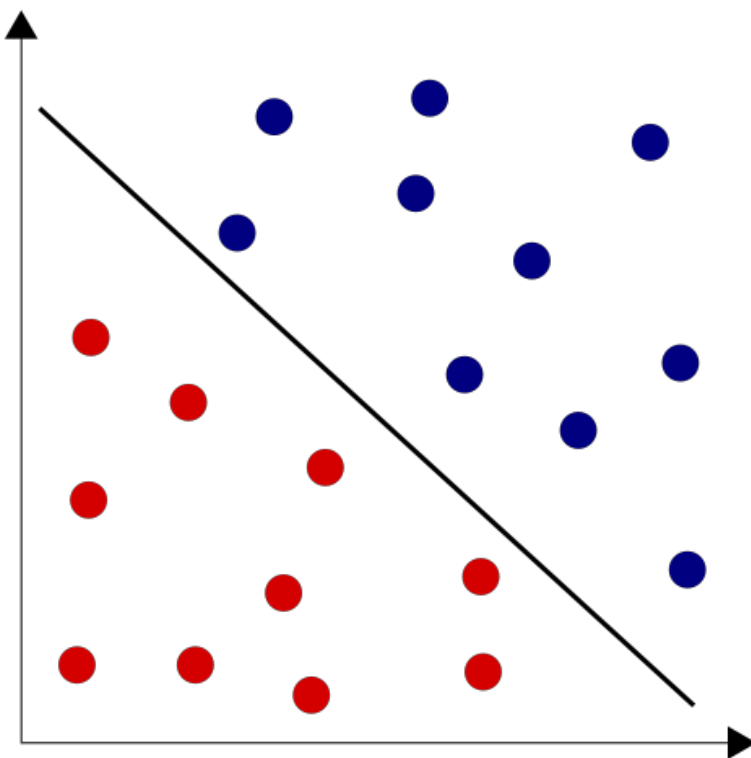
Власов Кирилл Вячеславович



Линейные модели в задаче классификации

Основная идея:

Предполагаем, что существует такая гиперплоскость, которая делит пространство на два полупространства в каждом из которых одно из двух значений целевого класса.



Если существует гиперплоскость которой можно разделить пространство на два класса без ошибок, то обучающая выборка называется *линейно разделимой*

Линейные модели в задаче классификации

Дана обучающая выборка:

$$X_1 = \{ (x_1, y_1), \dots, (x_1, y_1) \}$$

Для задачи классификации - Целевая переменная задана конечным числом меток

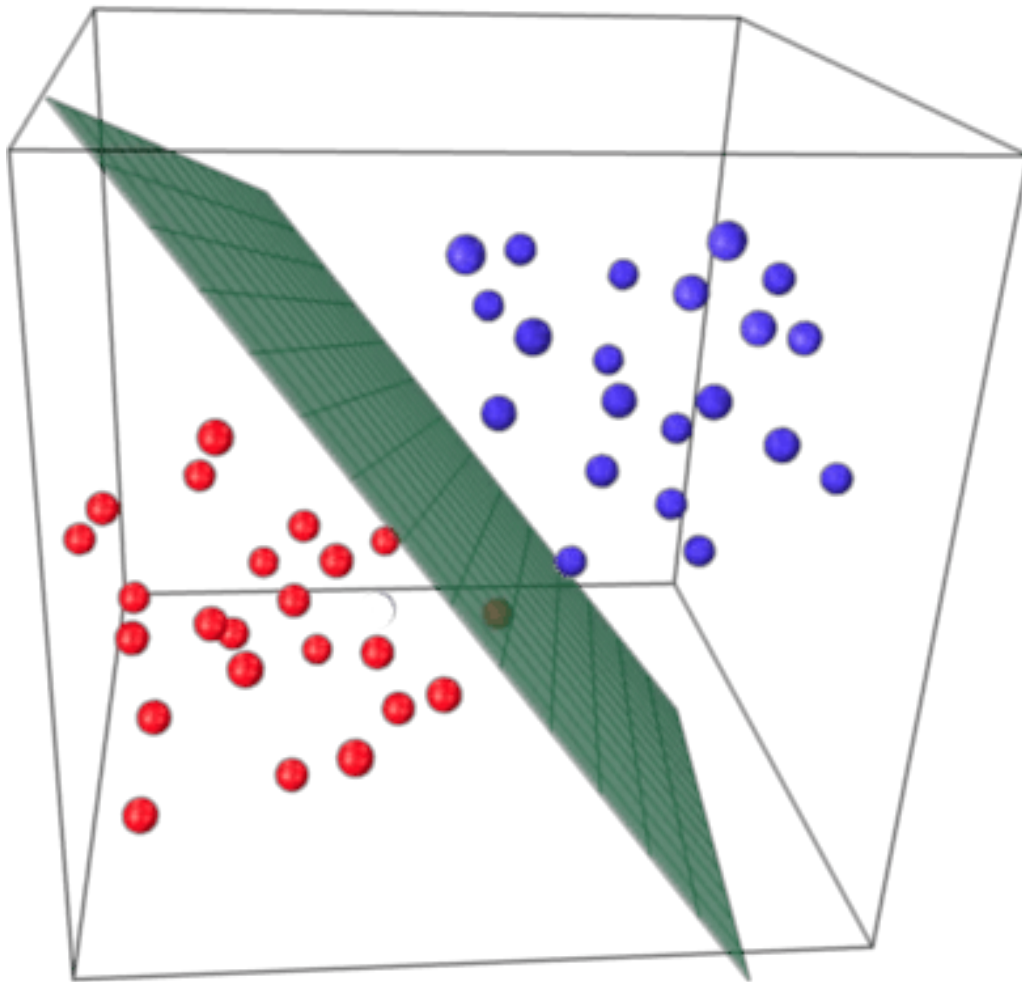
$$(x_1, y_1) \in \mathbb{R}^m \times \mathbb{Y}, \mathbb{Y} = \{-1; 1\}$$

Простейший классификатор:

$$a(x) = \text{sign}(\langle w, x \rangle + x_0) = \text{sign}(\vec{w}^T \cdot x)$$

\vec{w} – нормаль гиперплоскости

$\vec{w}^T \cdot x_i$ – расстояние от гиперплоскости до x_i ,
знак показывает отношение к классу



Рассмотрим пример!

Линейные модели в задаче классификации

$$a(x) = \text{sign}(\langle w, x \rangle + x_0) = \text{sign}(\vec{w}^T \cdot x)$$

доля правильных ответов (accuracy):

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) = y_i]$$

Линейные модели в задаче классификации

$$a(x) = \text{sign}(\langle w, x \rangle + x_0) = \text{sign}(\vec{w}^T \cdot x)$$

доля правильных ответов (accuracy):

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) = y_i] \rightarrow \max$$

Линейные модели в задаче классификации

$$a(x) = \text{sign}(\langle w, x \rangle + x_0) = \text{sign}(\vec{w}^T \cdot x)$$

доля правильных ответов (accuracy):

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) = y_i] \rightarrow \max$$

доля неправильных ответов:

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) \neq y_i] = \frac{1}{l} \sum_{i=1}^l [\text{sign}(\langle w, x_i \rangle) \neq y_i] \rightarrow \min$$

Проблемы:

1. Функционал дискретный относительно весов \Rightarrow мы не сможем искать минимум с помощью градиентных методов.
2. Функционал может иметь несколько глобальных минимумов \Rightarrow может быть много способов добиться оптимального количества ошибок.

Линейные модели в задаче классификации

$$a(x) = \text{sign}(\langle w, x \rangle + x_0) = \text{sign}(\vec{w}^T \cdot x)$$

доля правильных ответов (accuracy):

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) = y_i]$$

доля неправильных ответов:

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l [a(x_i) \neq y_i] = \frac{1}{l} \sum_{i=1}^l [\text{sign}(\langle w, x_i \rangle) \neq y_i] \rightarrow \min$$

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l \underbrace{[y_i \langle w, x_i \rangle < 0]}_{M_i} \rightarrow \min \quad M_i = y_i \langle w, x_i \rangle - \text{отступ (margin)}$$

*Знак отступа говорит о корректности ответа классификатора (положительный отступ соответствует правильному ответу, отрицательный неправильному)
абсолютная величина M – характеризует степень уверенности классификатора в своём ответе.*

Линейные модели в задаче классификации

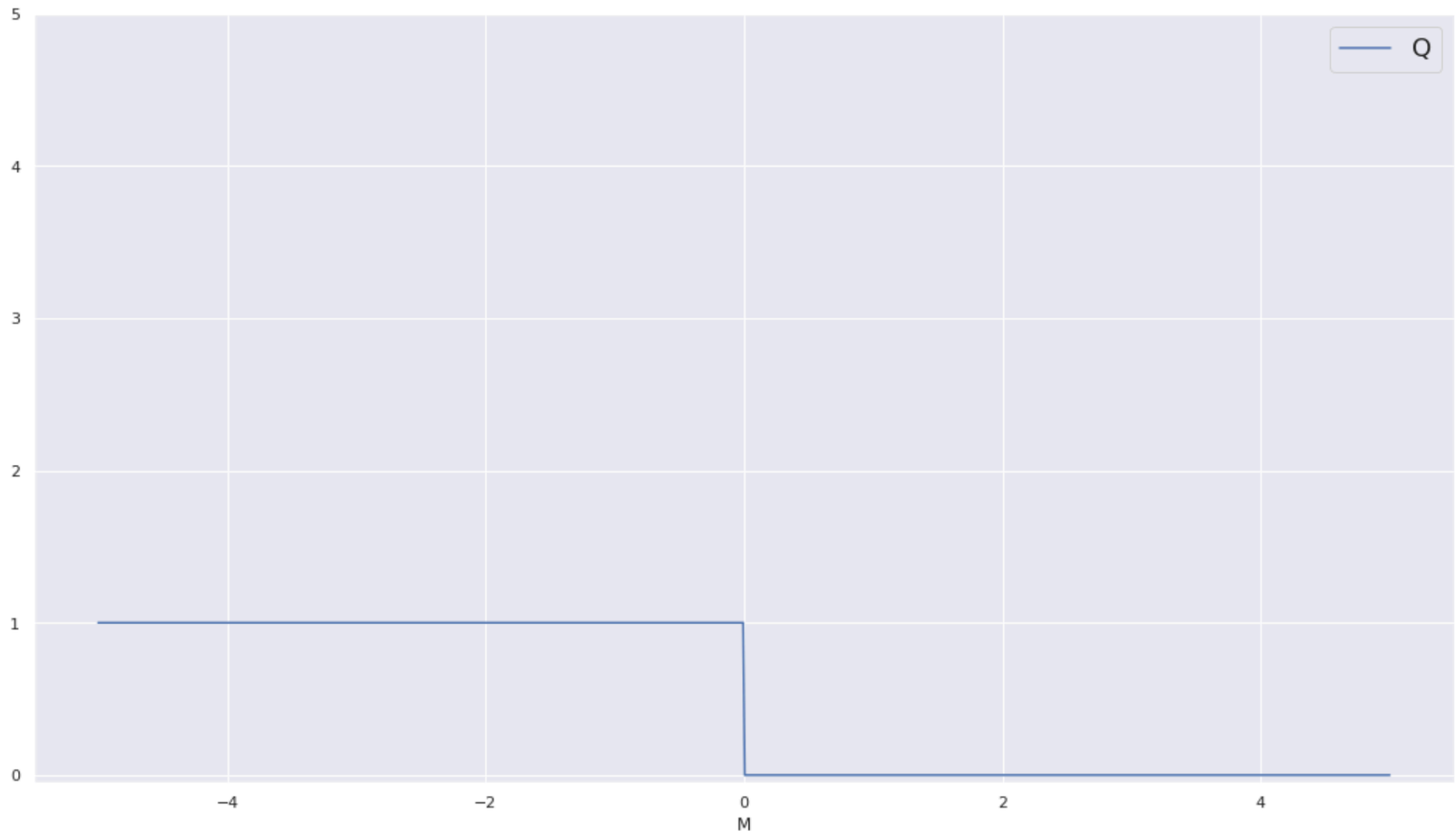
$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l \left[y_i \underbrace{\langle w, x_i \rangle}_{M_i} < 0 \right] \rightarrow \min$$

$$L(M) = [M < 0]$$

Линейные модели в задаче классификации

$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l \underbrace{[y_i \langle w, x_i \rangle < 0]}_{M_i} \rightarrow \min$$

$L(M) = [M < 0]$ – пороговая функции потерь



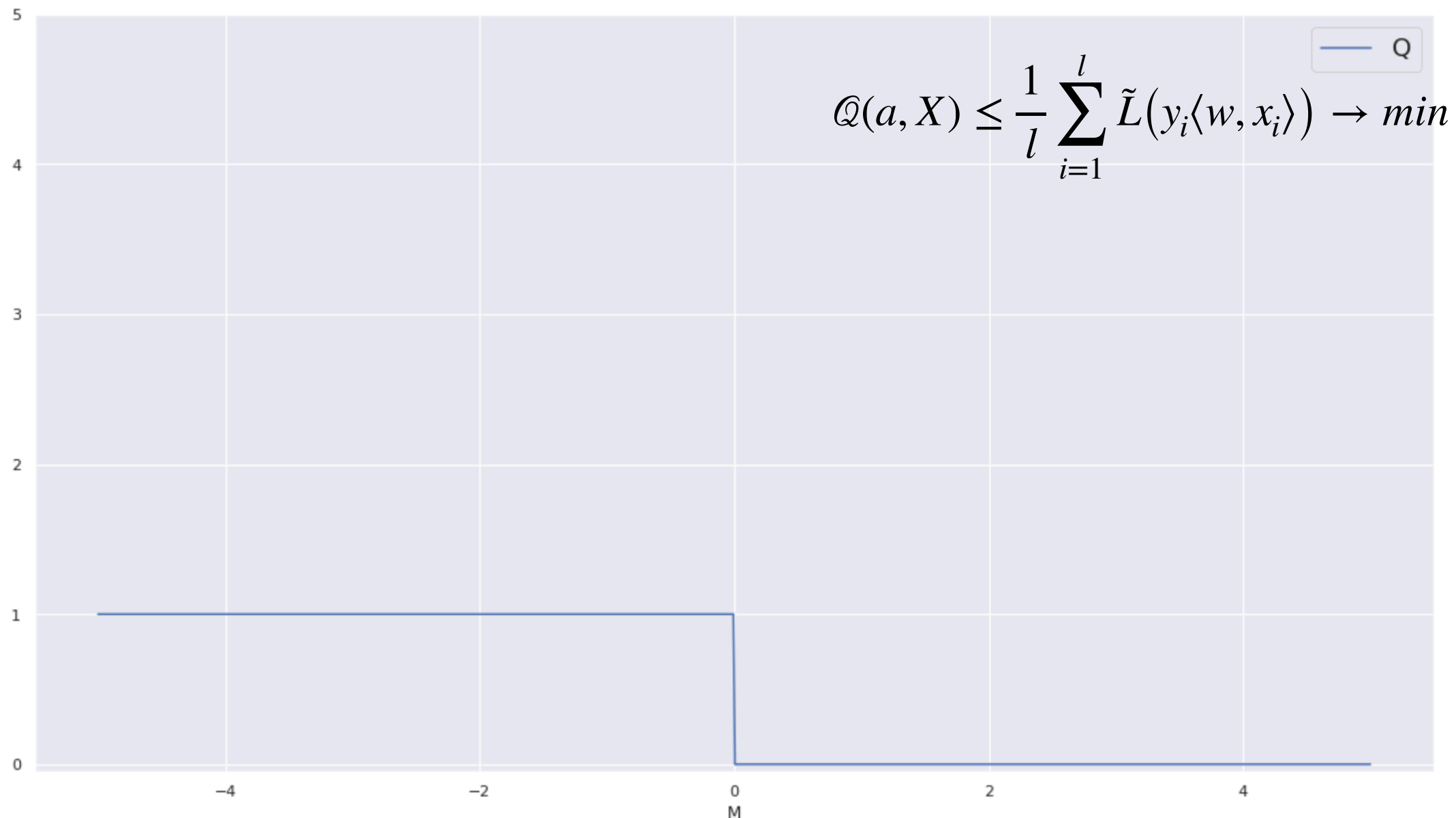
Линейные модели в задаче классификации

$$Q(a, X) = \frac{1}{l} \sum_{i=1}^l \underbrace{[y_i \langle w, x_i \rangle < 0]}_{M_i} \rightarrow \min$$

$L(M) \leq \tilde{L}(M)$ – верхняя оценка функции потерь

$L(M) = [M < 0]$ – пороговая функции потерь

Если верхнюю оценку удастся приблизить к нулю, то и доля неправильных ответов тоже будет близка к нулю



Линейные модели в задаче классификации

$$\tilde{L}(M) = (1 - M)^2$$

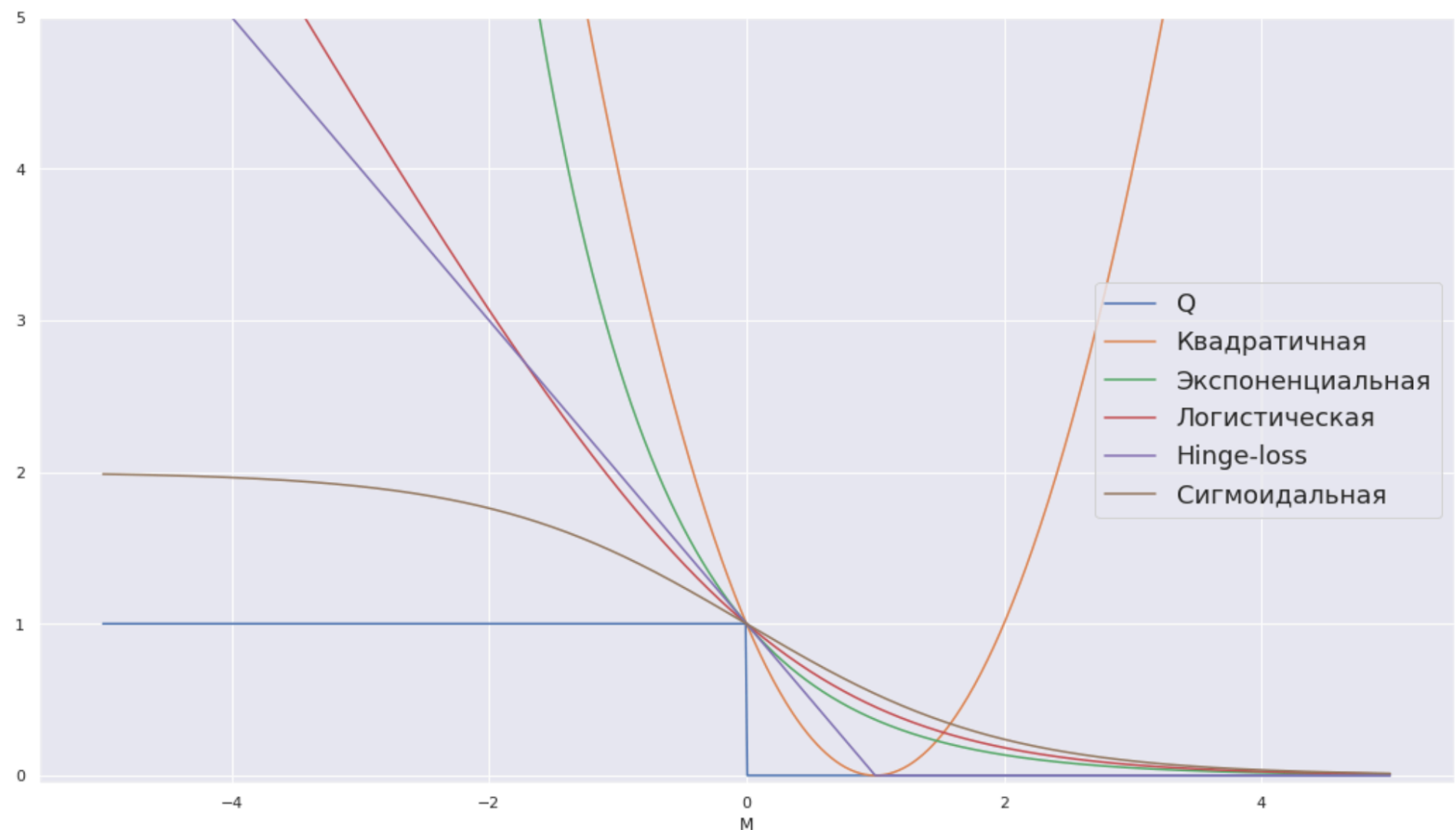
$$\tilde{L}(M) = e^{-M}$$

$$\tilde{L}(M) = \log(1 + e^{-M})$$

$$\tilde{L}(M) = (1 - M)_+ = \max(0, 1 - M)$$

$$\tilde{L}(M) = \frac{2}{1 + e^{-M}}$$

$L(M) = [M < 0]$ – пороговая функции потерь



Рассмотрим пример!

Линейные модели в задаче классификации

$$\tilde{L}(M) = (1 - M)^2$$

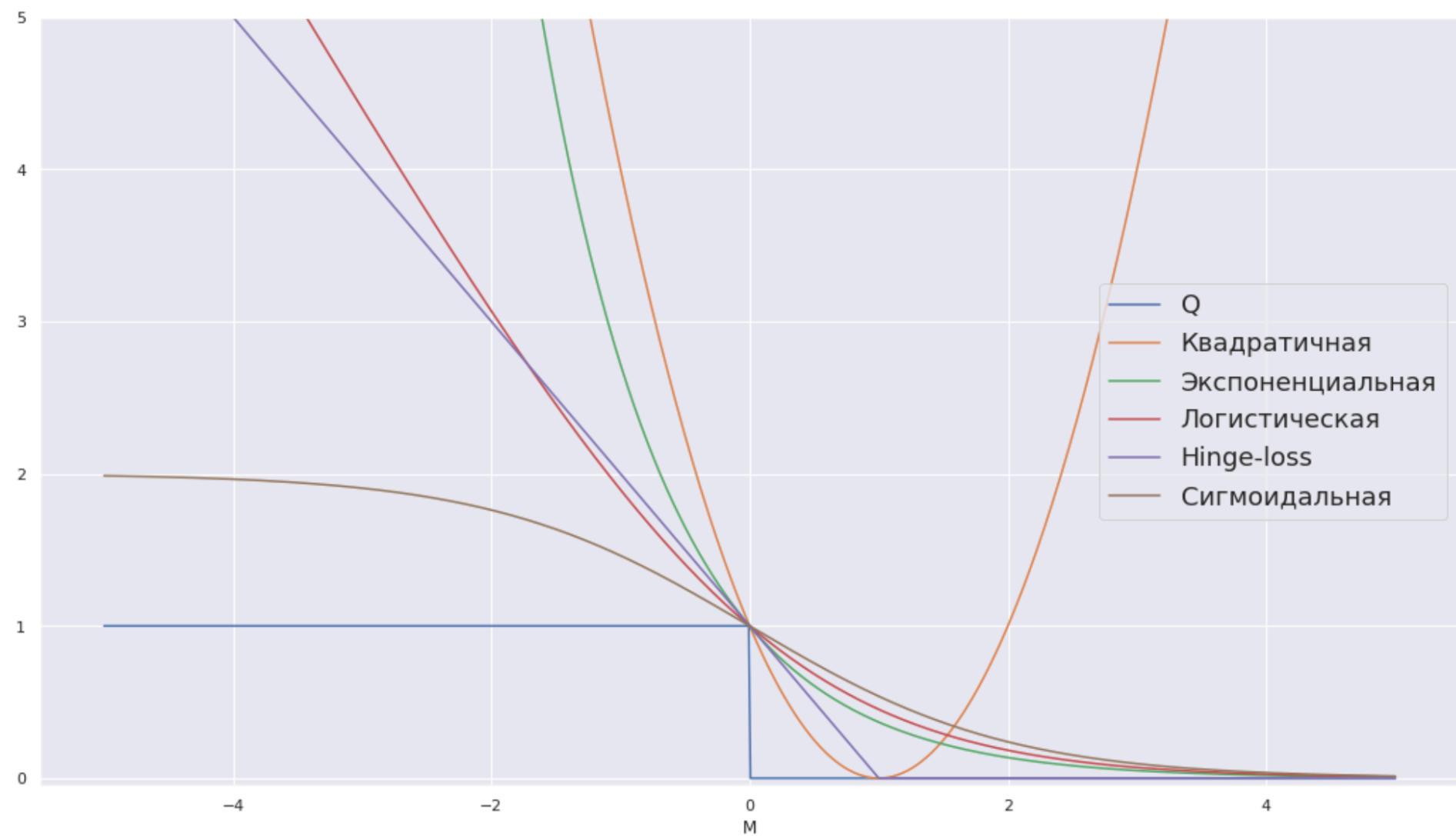
$$\tilde{L}(M) = e^{-M}$$

$$\tilde{L}(M) = \log(1 + e^{-M})$$

$$\tilde{L}(M) = (1 - M)_+ = \max(0, 1 - M)$$

$$\tilde{L}(M) = \frac{2}{1 + e^{-M}}$$

$L(M) = [M < 0]$ – пороговая функции потерь



Линейные модели в задаче классификации

$$\tilde{L}(M) = \log(1 + e^{-M})$$

Минимизация эмпирического риска:

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l \log(1 + e^{-y_i \langle w, x_i \rangle}) \rightarrow \min$$

$$a(x) = \text{sign}(\vec{w}^T \cdot x)$$

Линейные модели в задаче классификации

$$\tilde{L}(M) = \log(1 + e^{-M})$$

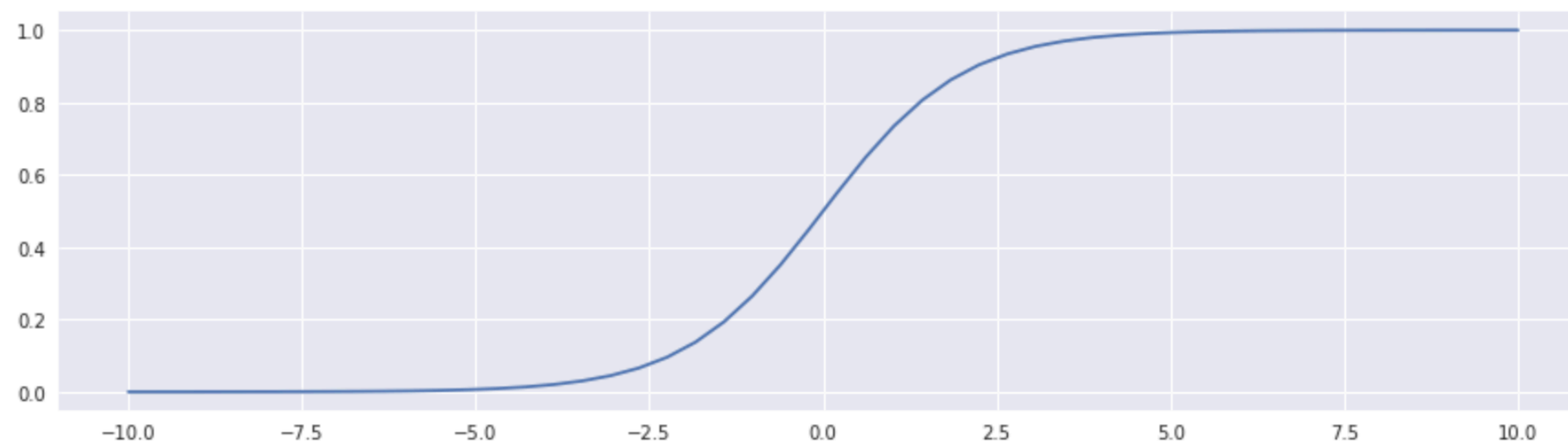
Минимизация эмпирического риска:

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l \log(1 + e^{-y_i \langle w, x_i \rangle}) \rightarrow \min$$

$$a(x) = \text{sign}(\vec{w}^T \cdot x)$$

Оценка апостериорных вероятностей принадлежности классам

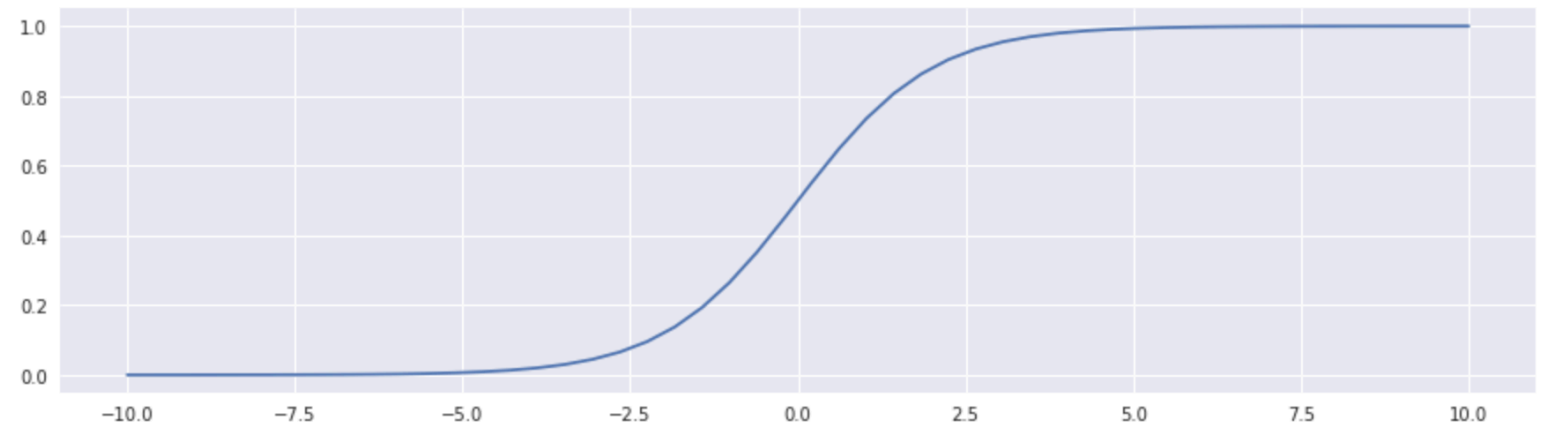
$$\vec{w}^T \cdot x \in \mathbb{R} \quad f(x_i, w) = \sigma(z) = \frac{1}{1 + e^{-z}} \in [0; 1]$$



Линейные модели в задаче классификации

$$y_i = f(x_i, w) + \varepsilon$$

$$f(x_i, w) = \sigma(z) = \frac{1}{1 + e^{-z}}$$



Шансы

$$odds = \frac{p}{1-p} \in [0; \infty] \quad \ln(odds) \in R$$

где p вероятность, что событие состоится (в нашем случае, что класс примет значение 1):

$$\ln(odds_+) = \ln\left(\frac{p}{1-p}\right) = \ln(p) - \ln(1-p)$$

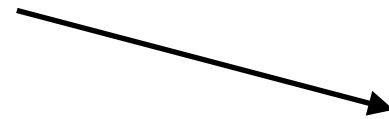
$$\ln(odds_-) = \ln\left(\frac{1-p}{p}\right) = \ln(1-p) - \ln(p)$$

$$\ln(odds_+) = -\ln(odds_-) = w^T \cdot x \quad odds = e^{w^T x} \Rightarrow p = \frac{e^{w^T x}}{1 + e^{w^T x}}$$

Линейные модели в задаче классификации

$$P(y_i = 1 | x_i, w) = \sigma(w^T x)$$

$$P(y_i = -1 | x_i, w) = \sigma(-w^T x)$$



$$P(y = y_i | x_i, w) = \sigma(y_i w^T x)$$

Правдоподобие (вероятность наблюдать вектор y при заданных значениях X и w)

Делаем предположение: объекты приходят независимо, из одного распределения

$$P(\vec{y} | X, w) = \prod_{i=1}^l P(y = y_i | x_i, w) \rightarrow \max$$

Так как логарифм монотонно возрастающая функция, то оценка w максимизирующая логарифм, будет максимизировать и само правдоподобие

$$\log P(\vec{y} | X, w) = \sum_{i=1}^l \log \sigma(y_i w^T x) = \sum_{i=1}^l \log \frac{1}{1 + e^{-y_i \langle w, x_i \rangle}} = - \sum_{i=1}^l \log(1 + e^{-y_i \langle w, x_i \rangle})$$

$$\mathcal{Q}(a, X) = \frac{1}{l} \sum_{i=1}^l \log(1 + e^{-y_i \langle w, x_i \rangle}) \rightarrow \min$$

Рассмотрим пример!

Ссылки

Открытый курс машинного обучения

Репозитории Евгения Соколова