

Обучение с учителем. Классификация. Функция риска. Логистическая регрессия. Feature selection и extraction

Третьякова Александра Леонидовна,
Волканова Маргарита Дмитриевна,
Федоров Никита Алексеевич



Санкт-Петербург
2019г.

X — множество объектов, Y — множество ответов
 $y : X \rightarrow Y$ — неизвестная зависимость (target function)

Дано: обучающая выборка — $(x_1, \dots, x_n) \subset X$,
 $y_i = y(x_i)$, $i = 1, \dots, n$ — известные ответы.
Найти: $a : X \rightarrow Y$ — функцию (decision function),
приближающую y на всем множестве X .

Вероятностная постановка задачи: имеется неизвестное распределение на множестве $X \times Y$ с плотностью $p(x, y)$, из которого случайно выбираются $\mathbb{X}_n = (x_i, y_i)_{i=1}^n$ (независимые).

Задача классификации:

- $Y = \{-1, +1\}$ — классификация на 2 класса
- $Y = \{1, \dots, K\}$ — классификация на K классов

Обучение с учителем. Постановка задачи

Обучающая выборка: $\mathbb{X}_n = (x_i, y_i)_{i=1}^n$, $x_i \in \mathbb{R}^p$, $y_i \in \{-1, +1\}$.

$f(x, w)$ — разделяющая (дискриминантная) функция, $w \in \mathbb{R}^p$.

$a(x, w) = \text{sign } f(x, w)$ — классификатор.

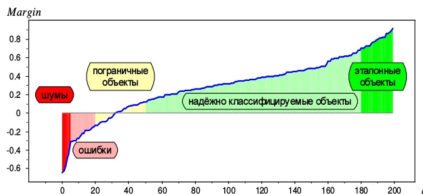
$f(x, w) = 0$ — разделяющая поверхность.

$M_i(w) = y_i f(x_i, w)$ — **отступ** объекта x_i .

Если $M_i(w) < 0$, то классификатор ошибается на x_i .

Задача

$$Q(w) = \sum_{i=1}^n [M_i(w) < 0] \leq \tilde{Q}(w) = \sum_{i=1}^n \mathcal{L}(M_i(w)) \rightarrow \min_w$$



Примеры функции потерь

- Пороговая функция потерь: $[M_i(w) < 0]$
- Логарифмическая: $\log_2(1 + e^{-M_i(w)})$
- Экспоненциальная: $e^{-M_i(w)}$
- Кусочно-линейная: $(1 - M_i(w))_+$
- Квадратичная: $(1 - M_i(w))^2$

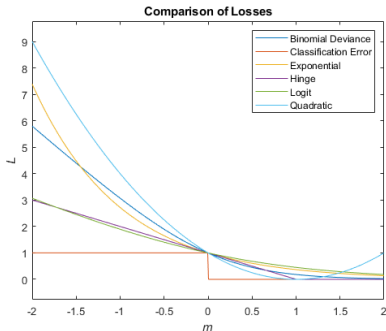


Рис.: Примеры функции потерь $\mathcal{L}(M)$

Линейный классификатор

$$f_j : X \rightarrow \mathbb{R}, j = 1, \dots, p$$

Линейная модель классификации:

$$a(x, w) = \text{sign}\left(\sum_{j=1}^p w_j f_j(x) - w_0\right), \quad w_0, w_1, \dots, w_p \in \mathbb{R}.$$

Пусть $f_0 = -1$, тогда

$$a(x, w) = \text{sign}\langle w, x \rangle, \quad x, w \in \mathbb{R}^{p+1}.$$

$M_i(x) = y_i \langle w, x_i \rangle$ — отступ объекта x_i .

Задача

$$Q(w) = \sum_{i=1}^n [y_i \langle w, x_i \rangle < 0] \leq \sum_{i=1}^n \mathcal{L}(y_i \langle w, x_i \rangle) \rightarrow \min_w$$

Проверка по тестовой выборке $\tilde{X}_k = (\tilde{x}_i, \tilde{y}_i)_{i=1}^k$:

$$\bar{Q}(w) = \frac{1}{k} \sum_{i=1}^k [\tilde{y}_i \langle w, \tilde{x}_i \rangle < 0]$$

Для решения задачи оптимизации используется
метод стохастического градиента.

Вход: $\mathbb{X}_n = (x_i, y_i)_{i=1}^n, h, \lambda$

Выход: w

- ❶ Инициализация $w_j, j = 0, \dots, p$
 - ❷ Инициализация $\bar{Q}(w)$
 - ❸ Повторять
 - ❶ Выбор x_i из \mathbb{X}_n случайным образом
 - ❷ Вычисление $\varepsilon_i := \mathcal{L}_i(w)$
 - ❸ Градиентный шаг $w := w - h \nabla \mathcal{L}_i(w)$
 - ❹ Вычисление $\bar{Q}(w) := \lambda \varepsilon_i + (1 - \lambda) \bar{Q}(w)$
- пока $\bar{Q}(w)$ или w не сойдутся

Проблемы:

- Признаков намного больше, чем объектов
- Мультиколлинеарность признаков:

Пусть

$$a(x, w) = \text{sign}(w_1 f_1(x) + w_2 f_2(x) - w_0), \quad f_2(x) = k f_1(x).$$

Тогда

$$w_1 f_1(x) + w_2 f_2(x) = (w_1 + \beta) f_1(x) + (w_2 - k\beta) f_2(x) \quad \forall \beta$$

Таким образом, очень много различных векторов дадут близкие значения функционала качества, но при этом коэффициенты могут существенно отличаться. Признаком такого явления может являться большая $\|w\|$.

Задача с регуляризацией

$$Q(w, \mathbb{X}_n) = Q(w, \mathbb{X}_n) + \frac{\tau}{2} \|w\|^2 \rightarrow \min_w$$

Связь с принципом максимума правдоподобия

Рассмотрим **модель**: пусть $X \times Y$ — вероятностное пространство с плотностью $p(x, y|w)$.

Пусть $\mathbb{X}_n = (x_i, y_i)_{i=1}^n \sim p(x, y|w)$ — независимы, одинаково распределены.

- **Максимизация правдоподобия:**

$$L(w; \mathbb{X}_n) = \ln \prod_{i=1}^n p(x_i, y_i|w) = \sum_{i=1}^n \ln p(x_i, y_i|w) \rightarrow \max_w.$$

- **Минимизация аппроксимированного эмпирического риска:**

$$\tilde{Q}(w; \mathbb{X}_n) = \sum_{i=1}^n \mathcal{L}(y_i f(x_i, w)) \rightarrow \min_w.$$

Эти задачи эквивалентны, если положить $-\ln p(x_i, y_i|w) = \mathcal{L}(y_i f(x_i, w))$.

Пример: $p(x, y|w) = \frac{1}{1 + \exp(-y\langle x, w \rangle)}$ (сигмоидная функция),
 $\mathcal{L}(yf(x, w)) = \log(1 + \exp(-y\langle x, w \rangle))$ (логарифмическая функция потерь)

Связь с принципом максимума правдоподобия

Модель: $\mathbb{X}_n = (x_i, y_i)_{i=1}^n \sim p(x, y|w)$ — н.о.р.,

Пусть $w \sim p(w; \gamma)$, γ — вектор гиперпараметров. Тогда:

$$p(\mathbb{X}_n, w) = p(\mathbb{X}_n|w)p(w; \gamma)$$

Принцип максимума правдоподобия:

$$L(w; \mathbb{X}_n) = \ln p(\mathbb{X}_n, w) = \sum_{i=1}^n p(x_i, y_i|w) + \underbrace{\ln p(w; \gamma)}_{\text{регуляризатор}} \rightarrow \max_{w, \gamma}$$

Примеры:

❶ **Гауссовский регуляризатор:** $p(w; \sigma) = \frac{1}{(2\pi\sigma)^{p/2}} \exp -\frac{\|w\|^2}{2\sigma}$,

тогда $-\ln p(w; \sigma) = \frac{1}{2\sigma} \|w\|^2 + \text{const}$ ($\tau = 1/\sigma$)

❷ **Регуляризатор Лапласа** (приводит к отбору признаков):

$$p(w; C) = \frac{1}{(2C)^p} \exp -\frac{\|w\|_1}{C},$$

тогда $-\ln p(w; C) = \frac{1}{C} \sum_{j=1}^p |w_j| + \text{const}$ ($\tau = 1/C$)

$$\text{Задача: } Q(w; \mathbb{X}_n) = \sum_{i=1}^n \ln p(x_i, y_i | w) + \frac{1}{C} \sum_{j=1}^p |w_j| \rightarrow \min_{w, C}.$$

$$\text{Замена: } \begin{cases} u_j = \frac{1}{2}(|w_j| + w_j) \\ v_j = \frac{1}{2}(|w_j| - w_j) \end{cases}, \text{ тогда } \begin{cases} w_j = u_j - v_j \\ |w_j| = u_j + v_j \end{cases},$$

$$\begin{cases} Q(u, v) = \sum_{i=1}^n \mathcal{L}(M_i(u - v, w_0)) + \frac{1}{C} \sum_{j=1}^p (u_j + v_j) \rightarrow \min_{u, v} \\ u_j \geq 0, v_j \geq 0, j = 1, \dots, p \end{cases}$$

При уменьшении C (возрастании $\frac{1}{C}$) обнуляются u_j и v_j для все большего количества j , то есть $w_j = 0$ и признак не учитывается.

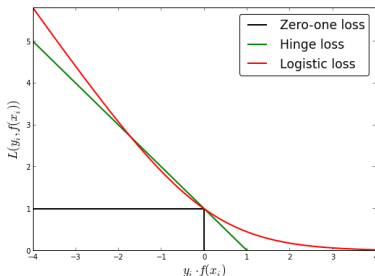
При $C \rightarrow 0$ выбросим все признаки.

Логистическая регрессия. Подход через минимизацию функции потерь

Линейная модель классификации:

$a(x) = \text{sign}\langle w, x \rangle$, $x, w \in \mathbb{R}^p$, $M = \langle w, x \rangle y$ — отступ.

В качестве аппроксимации пороговой функции потерь берется логарифмическая функция потерь $\mathcal{L}(M) = \log(1 + e^{-M})$.



Логистическая регрессия. Подход через минимизацию функции потерь

Задача

$$Q(w) = \sum_{i=1}^n \log(1 + \exp(-y_i \langle w, x_i \rangle)) \rightarrow \min_w$$

Методы решения задачи минимизации:

- метод стохастического градиента
- метод Ньютона-Рафсона

$P(y|x, w) = \sigma_w(M) = \frac{1}{1+e^{-\langle x, w \rangle y}}$ — сигмоидная функция.

Свойства $\sigma(z)$:

- $\sigma(z) \in [0, 1]$, задана на $(-\infty, +\infty)$
- $\sigma(z) \rightarrow 1, z \rightarrow +\infty;$
 $\sigma(z) \rightarrow 0, z \rightarrow -\infty$
- $\sigma(z) + \sigma(-z) = 1$
- $\sigma'(z) = \sigma(z)\sigma(-z)$

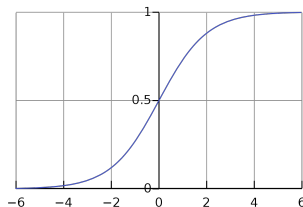


Рис.: Сигмоидная функция

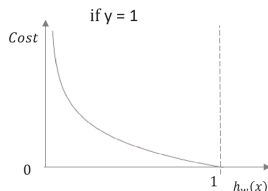
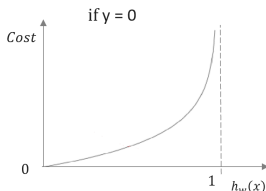
Пусть $Y = \{0, 1\}$.

- $P(y_i = 1|x; w) = \sigma_w(x)$
- $P(y_i = 0|x; w) = 1 - \sigma_w(x)$

Тогда $P(y|x; w) = (\sigma_w(x))^y(1 - \sigma_w(x))^{1-y}$.

Функция правдоподобия:

$$\begin{aligned} Q(w) &= -\log L(w) = -\log \prod_{i=1}^n (\sigma_w(x_i))^{y_i} (1 - \sigma_w(x_i))^{1-y_i} = \\ &= -\sum_{i=1}^n [y_i \log(\sigma_w(x_i)) + (1 - y_i) \log(1 - \sigma_w(x_i))] \rightarrow \min_w \end{aligned}$$



Линейная и логистическая регрессия

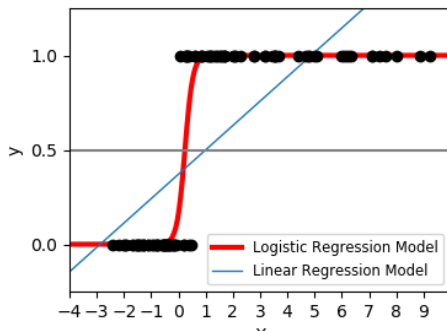


Рис.: Линейная и логистическая регрессия

$$Q(w) = - \sum_{i=1}^n [y_i \log(\sigma_w(x_i) + (1 - y_i)) \log(1 - \sigma_w(x_i))]$$

Регуляризация в логистической регрессии:

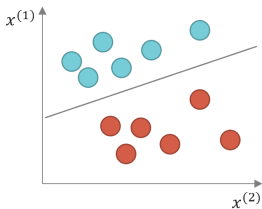
- L2: $Q_\tau(w) = Q(w) + \frac{\tau}{2} \sum_{j=1}^p w_j^2 \rightarrow \min_w$
- L1: $Q_\tau(w) = Q(w) + \tau \sum_{j=1}^p |w_j| \rightarrow \min_w$

Параметр τ можно подбирать с помощью кросс-валидации.

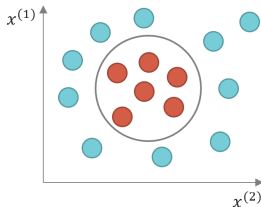
Методы решения задачи минимизации: метод стохастического градиента, метод Ньютона-Рафсона.

Логистическая регрессия. Добавление нелинейных признаков

Можно ли использовать логистическую регрессию в случае, когда нет линейной разделимости?



$\langle w, x \rangle = w_1 x_1 + w_2 x_2 + w_0$
Разделяющая поверхность:
 $w_1 x_1 + w_2 x_2 + w_0 = 0$



$\langle w, x \rangle = w_1 x_1^2 + w_2 x_2^2 + w_0$
Разделяющая поверхность:
 $w_1 x_1^2 + w_2 x_2^2 + w_0 = 0$

Линейный классификатор при произвольном числе классов $Y = \{1, \dots, K\}$:

$$a(x, w) = \arg \max_{y \in Y} \langle w_y, x \rangle, \quad x, w_y \in \mathbb{R}^p$$

Вероятность того, что объект x относится к классу i :

$$P(y = i | x; w) = \frac{\exp \langle w_i, x \rangle}{\sum_{z \in Y} \exp \langle w_z, x \rangle} = \frac{e^{w_i^T x}}{\sum_{k=1}^K e^{w_k^T x}}$$

Задача:

$$Q(w) = - \sum_{i=1}^n \log P(y_i | x_i; w) \rightarrow \min_w$$

Плюсы:

- ➊ Позволяет оценить вероятности принадлежности объектов к классу
- ➋ Достаточно быстро работает при больших объемах выборки
- ➌ Применима в случае отсутствия линейной разделимости, если на вход подать полиномиальные признаки

Минусы:

- ➊ Плохо работает в задачах, в которых зависимость сложная, нелинейная

Пример использования логистической регрессии. Задача кредитного скоринга

Пусть $Y = \{+1, -1\}$

Величина потери $D_{xy} = \begin{cases} S(x), & Y = -1 \text{ (кредит не вернули)} \\ -rS(x), & Y = +1 \text{ (кредит вернули)} \end{cases}$.

Логистическая регрессия дает возможность вычислять апостериорные вероятности принадлежности классу для каждого объекта x :

$$P(y|x; w) = \frac{1}{1 + e^{-\langle x, w \rangle y}}$$

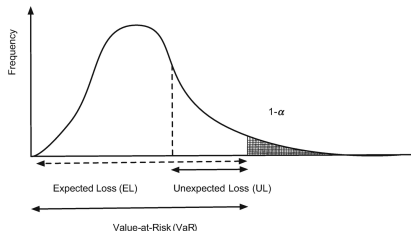
$R(x) = \sum_{y \in Y} D_{xy} P(y|x) = \sum_{y \in Y} D_{xy} \sigma_w(x)$ — оценка мат. ожидания потерь для объекта x . Хотим узнать, сколько банк потеряет в худшем случае.

Пример использования логистической регрессии. Задача кредитного скоринга

Строим эмпирическую функцию распределения потерь.

Метод Value at Risk:

- ❶ N раз ($N = 1000$):
 - $\forall x_i$ случайно разыгрываем $y_i \sim P(y|x_i)$, $i = 1, \dots, n$
 - вычисляем суммарные потери $V = \sum_{i=1}^n D_{x_i y_i}$
- ❷ строим эмпирическое распределение величины V
- ❸ 99%-квантиль показывает величину резервируемого капитала



- Существуют различные варианты аппроксимации пороговой функции потерь, позволяющие использовать методы градиентной оптимизации
- Регуляризация решает проблему мультиколлинеарности
- Минимизация аппроксимированного эмпирического риска и максимизация правдоподобия оказываются эквивалентными задачами
- Логистическая регрессия позволяет оценить условные вероятности классов
- В случае отсутствия линейной разделимости можно добавить нелинейные признаки и использовать логистическую регрессию