

Базовые модели машинного обучения: логистическая регрессия

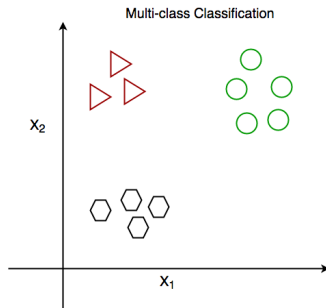
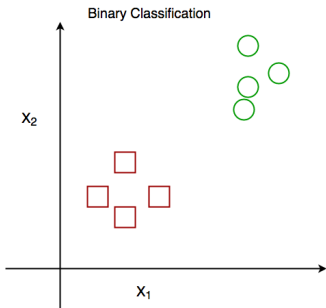
Першин Антон Юрьевич, Ph.D.
Никольская Анастасия Николаевна

Программа «Большие данные и распределенная цифровая
платформа»
Санкт-Петербургский государственный университет

Практика по дисциплине «Технологии ИИ»
31 марта 2023 г.

Задача классификации

- Пусть наблюдение характеризуется признаками $\mathbf{x} \in \mathbb{R}^M$
- Известно, что каждому наблюдению соответствует класс $y \in \mathcal{Y}$, где без потери общности $\mathcal{Y} = \{1, 2, \dots, K\}$
- Мы хотим построить правило $h : \mathbb{R}^M \rightarrow \mathcal{Y}$, позволяющее по наблюдению предсказывать его класс
- Предполагается, что имеется тренировочный набор данных, то есть набор пар $\mathcal{D} = \{(\mathbf{x}^{(i)}, y_i) | \mathbf{x}^{(i)} \in \mathbb{R}^M, y_i \in \mathcal{Y}\}_{i=1}^N$

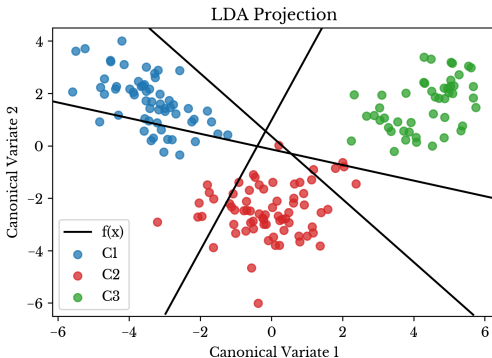


Задача классификации

- Если представить, что мы знаем апостериорное распределение $P(Y = y|x)$, то не составит труда задать правило $h(x)$ для предсказания:

$$h(x) = \operatorname{argmax}_{y \in Y} P(Y = y|x). \quad (1)$$

- Распределения неизвестны \implies можно оценить с помощью статистических методов
- Для пространств большой размерности их оценка затруднена, поэтому чаще используют подход с построением сепарирующей поверхности
- Например, линейный дискриминантный анализ использует линейную функцию для разделения классов



- Логистическая регрессия появляется из желания смоделировать апостериорное распределение линейной функцией, при этом гарантировав, что все вероятности суммируются в единицу и лежат в интервале $[0; 1]$
- Это возможно при использовании logit-преобразования:

$$\begin{aligned}\log \frac{P(Y = 1|X = \mathbf{x})}{P(Y = K|X = \mathbf{x})} &= \beta_{10} + \beta_1^T \mathbf{x}, \\ \log \frac{P(Y = 2|X = \mathbf{x})}{P(Y = K|X = \mathbf{x})} &= \beta_{20} + \beta_2^T \mathbf{x}, \\ &\dots \\ \log \frac{P(Y = K - 1|X = \mathbf{x})}{P(Y = K|X = \mathbf{x})} &= \beta_{K-1,0} + \beta_K^T \mathbf{x}.\end{aligned}$$

→ Так как сумма вероятностей должна быть равна единице, легко вывести функцию вероятности:

$$P(Y = y|X = \mathbf{x}) = \frac{\exp(\beta_{k0} + \beta_k^T \mathbf{x})}{1 + \sum_{j=1}^{K-1} \exp(\beta_{j0} + \beta_j^T \mathbf{x})}, \quad k = 1, \dots, K-1,$$

$$P(Y = K|X = \mathbf{x}) = \frac{1}{1 + \sum_{j=1}^{K-1} \exp(\beta_{j0} + \beta_j^T \mathbf{x})}.$$

→ Таким образом, распределение $P(Y = y|X = \mathbf{x}) = p_y(\mathbf{x}; \boldsymbol{\theta})$ параметризовано $\boldsymbol{\theta} = [\beta_{10}, \beta_1^T, \dots, \beta_{K-1,0}, \beta_{K-1}^T]$

→ Логарифмическая функция правдоподобия в таком случае имеет вид

$$l(\boldsymbol{\theta}) = \sum_{i=1}^N \log P(Y = y_i|X = \mathbf{x}^{(i)}) = \sum_{i=1}^N \log p_{y_i}(\mathbf{x}^{(i)}; \boldsymbol{\theta}). \quad (2)$$

Логистическая регрессия для бинарного случая

- В случае бинарной классификации нам достаточно одной функции для задания распределения. Для удобства заменим класс 2 на класс 0. Тогда

$$p(x; \theta) = P(Y = 1 | X = x) = \frac{1}{1 + \exp(\beta_0 + \beta^T x)} \quad (3)$$

так как $P(Y = 0 | X = x) = 1 - p(x; \theta)$. Здесь $\theta = [\beta_0, \beta]^T$

- Дополним x единицей в начале (смещение). Тогда логарифмическая функция правдоподобия может быть записана как

$$\begin{aligned} l(\theta) &= \sum_{i=1}^N \left[y_i \log p(x^{(i)}; \theta) + (1 - y_i) \log(1 - p(x^{(i)}; \theta)) \right] \\ &= \sum_{i=1}^N \left[y_i \beta^T x^{(i)} - \log(1 + \exp(\beta^T x^{(i)})) \right]. \end{aligned} \quad (4)$$

- Эта целевая функция выпуклая (\implies гарантирована сходимость к глобальному максимуму), но нелинейная (требуется метод Ньютона). Частичное доказательство выпуклости: (1) $\beta^T x^{(i)}$ вогнутая/выпуклая, \exp вогнутая и монотонно возрастающая $\implies 1 + \exp(\beta^T x^{(i)})$ вогнутая; LogSumExp вогнутая, тогда $-\text{LogSumExp}$ выпуклая.