

Санкт-Петербургский государственный университет

Прикладная математика и информатика

Статистическое моделирование

ОБУЧЕНИЕ С УЧИТЕЛЕМ. КЛАССИФИКАЦИЯ. ФУНКЦИЯ РИСКА.
ЛОГИСТИЧЕСКАЯ РЕГРЕССИЯ. FEATURE SELECTION И EXTRACTION

Третьякова Александра Леонидовна

Волканова Маргарита Дмитриевна

Федоров Никита Алексеевич

Санкт-Петербург

2020

Оглавление

1. Обучение с учителем. Постановка задачи

Пусть наблюдается некоторый количественный отклик Y и p предикторов (признаков) X_1, \dots, X_p . Будем предполагать, что между Y и $X = (X_1, \dots, X_p)$ существует определенная связь, которую можно представить в виде

$$Y = f^*(X) + \varepsilon,$$

где f^* — фиксированная, но неизвестная функция от предикторов, ε — ошибка, которая не зависит от X и имеет нулевое среднее значение.

Так как ошибки имеют нулевое среднее, можем предсказывать Y в соответствии с формулой

$$\hat{Y} = \hat{f}(X, \theta),$$

где \hat{f} — оценка f^* , θ — вектор параметров оценки, \hat{Y} — предсказанное значение Y .

Когда целью является предсказание, нам не важна точная форма функции \hat{f} , если она обеспечивает точные предсказания Y .

Пусть имеется выборка из n отдельных наблюдений: x_1, \dots, x_n для которых известны соответствующие значения отклика. Обозначим x_{ij} — значение j -го признака i -го наблюдения, $x_i = (x_{i1}, \dots, x_{ip})^\top$, y_i — отклик у i -го наблюдения.

Хотим найти такую функцию \hat{f} , что $y \approx \hat{f}(x)$ для любого наблюдения (x, y) . Совокупность X_n пар $(x_1, y_1), \dots, (x_n, y_n)$, которая участвует в оценке функции f^* , называется обучающей выборкой. Выборка $X'_k = (x'_i, y'_i)_{i=1}^k$, не участвующая в оценке функции f^* , называется тестовой (или контрольной).

Вероятностная постановка задачи: $Y \sim \text{Ber}(\sigma(X))$ при классификации на два класса и $Y \sim \text{Mult}(\sigma(X))$, где $\sigma(X) = (\sigma_1(X), \dots, \sigma_K(X))$ при классификации на K классов. Задача — построить оценку $\hat{\sigma} = \hat{\sigma}(X, \theta)$ параметра σ . Предсказание $\hat{Y} = \operatorname{argmax}_k \sigma_k(X)$.

2. Функция риска

Для оценки качества предсказания необходимо ввести функционал качества. Для этого сначала введём функцию потерь.

Определение 1. *Функция потерь (loss function) $L((x, y), \theta)$ — неотрицательная функция, характеризующая величину ошибки предсказания на объекте x .*

Для задачи классификации в качестве функции потерь обычно используется индикатор ошибки $L((x, y), \theta) = \delta(y, \hat{f}(x, \theta))$.

Теперь определим эмпирический риск — функционал, который используется для оценки качества предсказания.

Определение 2. Эмпирический риск $Q(X_n, \theta) = \frac{1}{n} \sum_{i=1}^n L((x_i, y_i), \theta)$

3. Регуляризация

Проблемы:

- Признаков намного больше, чем объектов

- Мультиколлинеарность признаков:

Пусть $\hat{f}(x, \theta) = \text{sign}(\theta_1 f_1(x) + \theta_2 f_2(x) - \theta_0)$, $f_2(x) = k f_1(x)$.

Тогда $\theta_1 f_1(x) + \theta_2 f_2(x) = (\theta_1 + \beta) f_1(x) + (\theta_2 - k\beta) f_2(x) \quad \forall \beta$

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

Задача с регуляризацией:

$$\tilde{Q}(X_n, \theta) = Q(X_n, \theta) + \frac{\tau}{2} \|\theta\|^2 \rightarrow \min_{\theta}$$

4. Связь с принципом максимума правдоподобия

Рассмотрим **модель**: пусть Y — вероятностное пространство с плотностью $p(y|x, \theta)$, то есть $\mathcal{L}(Y)$ — семейство распределений, зависящее от параметров X и θ .

Пусть $X_n = (x_i, y_i)_{i=1}^n$, и y_i — независимы, одинаково распределены.

- **Максимизация правдоподобия:**

$$\mathcal{L}(\theta; X_n) = \ln \prod_{i=1}^n p(y_i|x_i, \theta) = \sum_{i=1}^n \ln p(y_i|x_i, \theta) \rightarrow \max_{\theta}.$$

- **Минимизация аппроксимированного эмпирического риска:**

$$\tilde{Q}(X_n, \theta) = \sum_{i=1}^n L(x_i, y_i, \theta) \rightarrow \min_{\theta}.$$

Эти задачи эквивалентны, если положить $-\ln p(y_i|x_i, \theta) = L(x_i, y_i, \theta)$.

Пример:

- $p(y|x, \theta) = \frac{1}{1+\exp(-y\langle x, \theta \rangle)}$ — сигмоидная функция.
- $L(x, y, \theta) = \log(1 + \exp(-y\langle x, \theta \rangle))$ — логарифмическая функция потерь.

Принцип максимума правдоподобия: Пусть $\theta \sim p(\theta; \gamma)$.

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

Примеры:

1. Гауссовский регуляризатор:

$$p(\theta; \sigma) = \frac{1}{(2\pi\sigma)^{p/2}} \exp -\frac{\|\theta\|^2}{2\sigma},$$

тогда

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

2. Регуляризатор Лапласа (приводит к отбору признаков):

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

тогда

$$-\ln p(\theta; C) = \frac{1}{C} \sum_{j=1}^p |\theta_j| + \text{const} \quad (\tau = 1/C)$$

5. Регуляризатор Лапласа. Отбор признаков

Задача:

$$Q(X_n, \theta) = \sum_{i=1}^n \ln p(y_i|x_i, \theta) + \frac{1}{C} \sum_{j=1}^p |\theta_j| \rightarrow \min_{\theta, C}.$$

Замена:

$$\begin{cases} u_j = \frac{1}{2}(|\theta_j| + \theta_j) \\ v_j = \frac{1}{2}(|\theta_j| - \theta_j) \end{cases},$$

тогда

$$\begin{cases} \theta_j = u_j - v_j \\ |\theta_j| = u_j + v_j \end{cases}$$

$$\begin{cases} Q(u, v) = \sum_{i=1}^n \mathcal{L}(M_i(u - v, \theta_0)) + \frac{1}{C} \sum_{j=1}^p (u_j + v_j) \rightarrow \min_{u, v} \\ u_j \geq 0, v_j \geq 0, j = 1, \dots, p \end{cases}$$

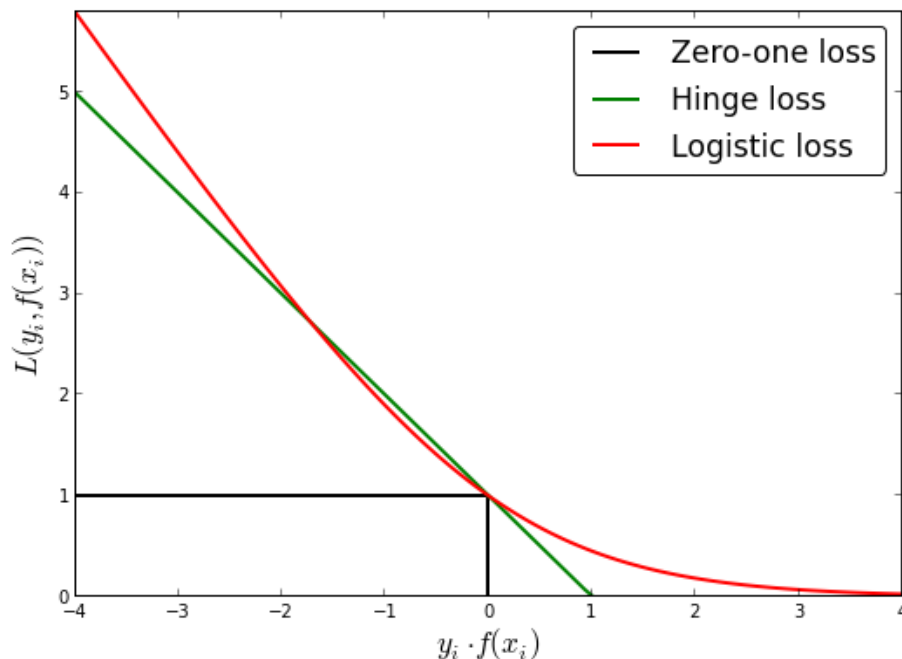
При уменьшении C (возрастании $\frac{1}{C}$) обнуляются u_j и v_j для все большего количества j , то есть $\theta_j = 0$ и признак не учитывается. При $C \rightarrow 0$ выбросим все признаки.

6. Логистическая регрессия. Подход через минимизацию функции потерь

Линейная модель классификации:

- $\hat{f}(x) = \text{sign}\langle \theta, x \rangle$, $x, \theta \in \mathbb{R}^p$
- $M = \langle \theta, x \rangle y$ — отступ.

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



Задача 1. $Q(X_n, \theta) = \sum_{i=1}^n \log(1 + \exp(-y_i \langle \theta, x_i \rangle)) \rightarrow \min_{\theta}$

Методы решения задачи минимизации:

- метод стохастического градиента
- метод Ньютона-Рафсона

7. Логистическая регрессия. Вероятностный подход

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

Свойства $\sigma(z)$:

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

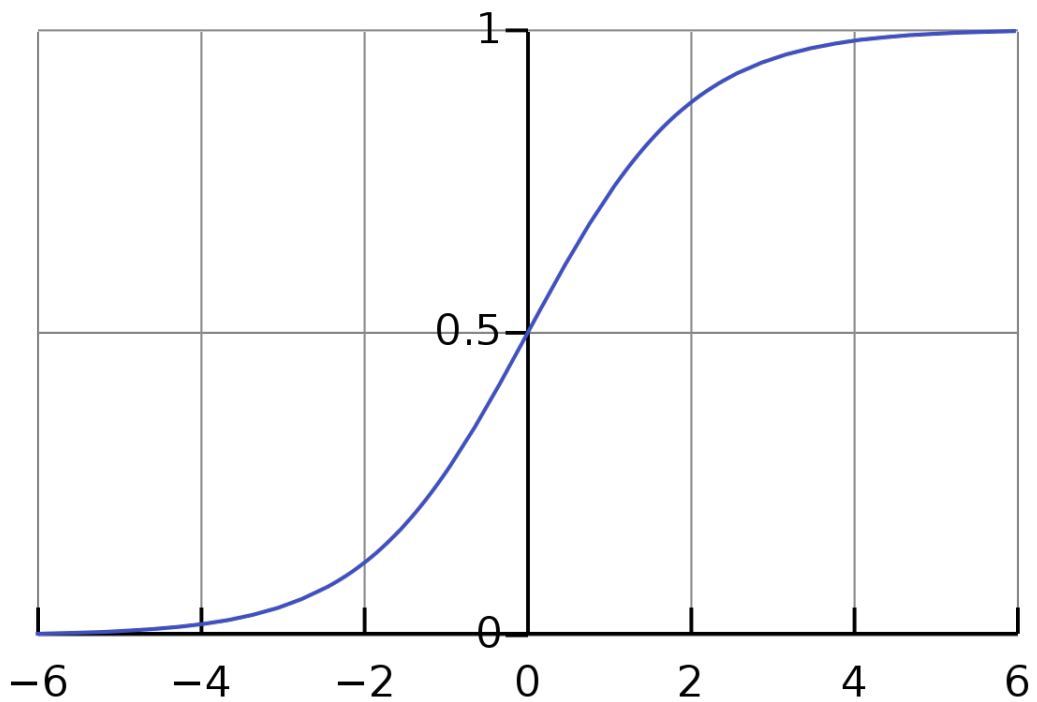


Рис. 1. Сигмоидная функция

Пусть $Y = \{0, 1\}$.

- $P(y_i = 1|x; \theta) = \sigma_{\theta}(x)$

- $P(y_i = 0|x; \theta) = 1 - \sigma_\theta(x)$

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

Функция правдоподобия:

$$\begin{aligned} Q(X_n, \theta) &= -\log L(\theta) = -\log \prod_{i=1}^n (\sigma_\theta(x_i))^{y_i} (1 - \sigma_\theta(x_i))^{1-y_i} = \\ &= -\sum_{i=1}^n [y_i \log(\sigma_\theta(x_i)) + (1 - y_i) \log(1 - \sigma_\theta(x_i))] \rightarrow \min_{\theta} \end{aligned}$$

8. Линейная и логистическая регрессия

Существуют примеры данных, для которых логистическая регрессия показывает лучшие результаты, чем линейная.

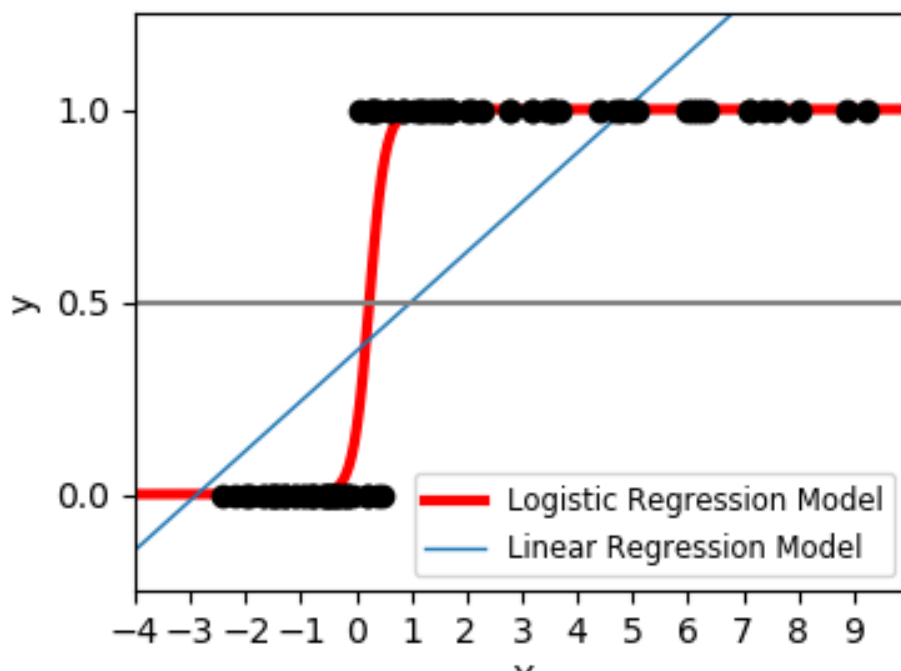


Рис. 2. Линейная и логистическая регрессия

9. Логистическая регрессия. Регуляризация

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

Регуляризация в логистической регрессии:

- **L2:** $Q_{\tau}(\theta) = Q(\theta) + \frac{\tau}{2} \sum_{j=1}^p \theta_j^2 \rightarrow \min_{\theta}$
- **L1:** $Q_{\tau}(\theta) = Q(\theta) + \tau \sum_{j=1}^p |\theta_j| \rightarrow \min_{\theta}$

Параметр τ можно подбирать с помощью кросс-валидации.

Методы решения задачи минимизации:

- метод стохастического градиента
- метод Ньютона-Рафсона.

10. Многоклассовая логистическая регрессия

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

$$\hat{f}(x, \theta) = \arg \max_{y \in Y} \langle \theta_y, x \rangle, \quad x, \theta_y \in \mathbb{R}^p$$

Вероятность того, что объект x относится к классу i :

$$P(y = i | x; \theta) = \frac{\exp \langle \theta_y, x \rangle}{\sum_{z \in Y} \exp \langle \theta_z, x \rangle} = \frac{e^{\theta_i^T x}}{\sum_{k=1}^K e^{\theta_k^T x}}$$

Задача:

$$Q(X_n, \theta) = - \sum_{i=1}^n \log P(y_i | x_i; \theta) \rightarrow \min_{\theta}$$

11. Логистическая регрессия. Преимущества и недостатки

Плюсы:

1. Позволяет оценить вероятности принадлежности объектов к классу
2. Достаточно быстро работает при больших объемах выборки
3. Применима в случае отсутствия линейной разделимости, если на вход подать полиномиальные признаки

Минусы:

1. Плохо работает в задачах, в которых зависимость сложная, нелинейная

12. Выводы

- 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