

Оценка оптимального объёма выборки для задач классификации*

Харатьян А. С., Катруца А. М.¹, Стрижов В. В.²

haratyan.as@phystech.edu; aleksandr.katrutsa@phystech.edu;
strijov@phystech.edu

¹Московский физико-технический институт, Москва, Россия; ²Вычислительный центр им. А. А. Дородницына ФИЦ ИУ РАН, Москва, Россия

В статье рассматривается задача выбора оптимального числа объектов выборки для их классификации. Исследуется использование порождающих и разделяющих вероятностных моделей бинарной классификации. Обсуждается проблема медицинской диагностики пациентов. Определяется понятие достаточности объёма выборки. Показывается, какими методами возможно выбрать оптимальное количество объектов, обеспечивающее необходимую точность классификации объектов. В работе рассматривается, применение каких критериев выявляет наилучшее качество классификации. Приводится теоретическое и практическое обоснование предложенных критериев. Используется модель логистической регрессии.

Ключевые слова: *определение оптимального объёма выборки, логистическая регрессия, расстояние Кульбака-Лейблера.*

1 Введение

Работа посвящена оценке оптимального объёма исследуемой выборки применительно к проблемам медицинской диагностики. Рассматриваются биомедицинские данные пациентов как выборка. Каждый пациент набором признаков. Получение данных о пациентах требует немалых средств. В случае, если количество данных избыточно, то их измерения приносят крайне неоправданные расходы. В связи с этим поднимается вопрос оптимального количества измерений. Ввиду дороговизны анализов всех признаков оценка измерений должна быть точной.

Для нахождения оценки используется модель логистической регрессии [1]. Стандартной практикой является использование статистических методов [2] для оценивания объёма данных при помощи логистической регрессии. Введём понятие устойчивости модели в отношении объёма выборки. Будем называть модель устойчивой, если при изменении малом изменении объёма параметры модели меняются незначительно. Если размер выборки крайне мал и недостаточен, то параметры модели меняются скачкообразно при увеличении объёма. Соответственно, увеличивая объём выборки, мы повышаем устойчивость модели. В качестве показателя устойчивости для моделей будем использовать расстояние Кульбака-Лейблера. Для того чтобы показать отличие моделей в устойчивости будем использовать разность усреднённых значений расстояний К-Л, вычисленных на разных выборках одного и того же объёма. Если объекты порождены одинаковым распределением, то при росте объёма выборки разность расстояний К-Л между моделями падает. Достигнув необходимого показателя устойчивости, можно легко вычислить оптимальный размер данных.

*Работа выполнена при финансовой поддержке РФФИ, проект №00-00-00000. Научный руководитель: Стрижов В. В. Консультант: Катруца А. М.

Для вычислительного эксперимента используются реальные данные 569 пациентов с 30 признаками и метками об опухоли молочной железы: доброкачественная или злокачественная. Как описано выше, будем обучать модели на разных подвыборках и после достижения необходимого уровня устойчивости получим оптимальный объём данных.

2 Постановка задачи классификации

Пусть у нас задана выборка $D = \{(\mathbf{x}_i, y_i) : i = 1, \dots, m\}$ с объёмом m объектов (пациентов), каждый из которых описывается n признаками, $\mathbf{x}_i \in \mathbb{R}^n$ и принадлежит одному из двух классов: $y_i \in \{0, 1\}$. Модель логистической регрессии предполагает, что вектор целевой переменной $\mathbf{y} = [y_1, \dots, y_m]^T$ имеет распределение Бернулли, $y_i \sim \mathcal{B}(\theta_i)$ с плотностью распределения

$$p(y|\boldsymbol{\omega}) = \prod_{i=1}^m \theta_i^{y_i} (1 - \theta_i)^{1-y_i} \quad (1)$$

Плотность вероятности зависит от вектора параметров $\boldsymbol{\omega}$. Зная $\boldsymbol{\omega}$, можно вычислить вероятность принадлежности к классу

$$\theta_i = f(\mathbf{x}_i^T \boldsymbol{\omega}) = \frac{1}{1 + \exp(-\mathbf{x}_i^T \boldsymbol{\omega})} \quad (2)$$

Пользуясь принципом максимального правдоподобия, можем вычислить функцию ошибки для уравнения (1)

$$E(\boldsymbol{\omega}) = -\ln p(\mathbf{y}|\boldsymbol{\omega}) = -\sum_{i=1}^m (y_i \ln \theta_i + (1 - y_i) \ln (1 - \theta_i)) \quad (3)$$

Чтобы найти вектор параметров логистической регрессии $\hat{\boldsymbol{\omega}}$, необходимо решить оптимизационную задачу:

$$\hat{\boldsymbol{\omega}} = \arg \min_{\boldsymbol{\omega} \in \mathbb{R}^n} E(\boldsymbol{\omega}) \quad (4)$$

Тогда алгоритм классификации определяется следующим образом:

$$a(\mathbf{x}, c_0) = \text{sign}(f(\mathbf{x}, \boldsymbol{\omega}) - c_0) \quad (5)$$

где c_0 - пороговое значение функции активации [3].

3 Постановка задачи определения оптимального объёма выборки

Допустим, что задано m объектов выборки D . Количество этих объектов недостаточно для обучения устойчивой модели логистической регрессии. Необходимо, имея m объектов, найти такое число m^* , что показатель устойчивости модели при увеличении числа m^* меняется незначительно. Задача нахождения данного числа объектов представляется следующим образом:

$$m^* = \min m \in \mathbb{N} \quad \forall k \in \mathbb{N} \rightarrow \rho(p(\mathbf{w}|m^*), p(\mathbf{w}|m^* + k)) < \varepsilon \quad (6)$$

где ρ - некоторая функция расстояния между распределениями, ε - заранее заданный порог устойчивости, $p(\mathbf{w}|m^*)$ - распределение весов модели, обученной на выборке размером m^* .

Литература

- [1] Hosmer Jr, D. W., S. Lemeshow, and R. X. Sturdivant. 2013. *Applied logistic regression*. John Wiley & Sons. Vol. 398.
- [2] Demidenko, E. 2007. Sample size determination for logistic regression revisited. *Statistics in medicine* 26(18):3385–3397.
- [3] Motrenko, A., V. Strijov, and G.-W. Weber. 2014. Sample size determination for logistic regression. *Journal of Computational and Applied Mathematics* 255:743–752.

Received