

Раннее прогнозирование достаточного объема выборки для обобщенной линейной модели.

Бучнев В. С., Грабовой А. В., Гадаев Т. Т., Стрижов В. В.

Исследуется проблема снижения затрат на сбор данных, необходимых для построения адекватной модели. Рассматриваются задачи линейной и логистической моделей. Для решения этих задач требуется, чтобы выборка содержала необходимое число объектов. Требуется предложить метод вычисления оптимального объема данных, соблюдая при этом баланс между точностью модели и и трудозатратами при сборе данных. Предпочтительны те методы оценки объема, которые позволяют строить адекватные модели по выборкам возможно меньшего объема.

Ключевые слова: *Обобщенная линейная модель, размер выборки.*

1 Введение

При планировании эксперимента требуется оценить минимальный объем выборки — число производимых измерений набора показателей или признаков, необходимый для построения сформулированных условий.

Существует большое количество оценки размера выборки. Например, тест множителей Лагранжа, тест отношения правдоподобия и тест Вальда. В работах [1–3] на основе данных методов построена оценка оптимального размера выборки. Основным минус этих методов заключается в том, что статистики, используемые в критериях, имеют асимптотическое распределение и требуют большого объема выборки.

Существуют также байесовские оценки объема выборки: критерий средней апостериорной дисперсии, критерий среднего покрытия, критерий средней длины и метод максимизации полезности. Первые три метода требуют анализа некоторой функции эффективности от размера выборки. Используя некоторое решающее правило, по данной функции определяется достаточный объем выборки. Главный минус этих методов заключается в том, что они не позволяют построить аппроксимацию функции эффективности при большем объеме данных. Метод максимизации полезности максимизирует ожидание некоторой функции полезности по объему выборки. Все эти методы опираются на апостериорное распределение, что требует достаточно большого объема выборки.

Предлагается исследовать зависимость среднего значения логарифма правдоподобия от размера доступной выборки, а также его дисперсию. В данной работе предлагается использовать не сами функции эффективности, а их аппроксимации. Для этого предлагается использовать аппроксимацию ковариационной матрицы вектора параметров. После чего аппроксимировать данные две зависимости при помощи метода бутстреп. Для вычислительного эксперимента предлагается использовать классические выборки из UCI репозитория и синтетические данные.

2 Постановка задачи

Дана выборка размера m :

$$\mathcal{D}_m = \{\mathbf{x}_i, y_i\}_{i=1}^m,$$

где $\mathbf{x}_i \in \mathbb{R}^n$ - вектор признаков, $y_i \in \mathbb{Y}$.

Выборка \mathcal{D}_m разбита случайно на обучение и контроль:

$$\mathcal{D}_{\mathcal{T}_m} = \{\mathbf{x}_i, y_i\}_{i \in \mathcal{T}_m} \quad \mathcal{D}_{\mathcal{L}_m} = \{\mathbf{x}_i, y_i\}_{i \in \mathcal{L}_m}, \quad \mathcal{T}_m \sqcup \mathcal{L}_m = \{1, \dots, m\}.$$

Предполагается, что выборка \mathfrak{D}_m не противоречит гипотезе порождения данных.

Рассмотрим параметрическое семейство функций для аппроксимации неизвестного распределения $p(y|\mathbf{x}, \mathbf{w})$, где $\mathbf{w} \in \mathbb{W}$ - вектор параметров:

$$\mathfrak{F} = \left\{ f(y, \mathbf{x}, \mathbf{w}) | \mathbf{w} \in \mathbb{W}, \int_{y \in \mathbb{Y}, \mathbf{x} \in \mathbb{R}^n} f(y, \mathbf{x}, \mathbf{w}) dy d\mathbf{x} = 1 \right\}.$$

Для модели f с вектором параметров \mathbf{w} определим функцию правдоподобия и логарифмическую функцию правдоподобия выборки \mathfrak{D} :

$$L(\mathfrak{D}, \mathbf{w}) = \prod_{y, \mathbf{x} \in \mathfrak{D}} f(y, \mathbf{x}, \mathbf{w}), \quad l(\mathfrak{D}, \mathbf{w}) = \sum_{y, \mathbf{x} \in \mathfrak{D}} \log f(y, \mathbf{x}, \mathbf{w}),$$

где $f(y, \mathbf{x}, \mathbf{w})$ - аппроксимация апостериорной вероятности выборки $\mathfrak{D}_{\mathcal{L}_m}$ при заданном векторе параметров \mathbf{w} .

Рассмотрим правдоподобие выборки $\mathfrak{D}_{\mathcal{L}_m}$:

$$L(\mathfrak{D}_{\mathcal{L}_m}, \mathfrak{D}_{\mathcal{L}_m}) = \prod_{y, \mathbf{x} \in \mathfrak{D}_{\mathcal{L}_m}} f(y, \mathbf{x}, \mathbf{w}).$$

Рассмотрим логарифм правдоподобия выборки $\mathfrak{D}_{\mathcal{L}_m}$:

$$l(\mathfrak{D}_{\mathcal{L}_m}, \mathbf{w}) = \sum_{y, \mathbf{x} \in \mathfrak{D}_{\mathcal{L}_m}} \log f(y, \mathbf{x}, \mathbf{w}).$$

Будем рассматривать ожидаемое значение функции l :

$$\tilde{l}(\mathfrak{D}) = \mathbb{E}_{y, \mathbf{x} \in \mathfrak{D}} l(\{y, \mathbf{x}\}, \mathbf{w}).$$

Рассмотрим ожидаемое значение логарифма правдоподобия по разным обучающим выборкам $\mathfrak{D}_{\mathcal{L}_m}$ размера m^* :

$$l(m^*) = \mathbb{E}_{\mathfrak{D}_{\mathcal{L}_m}} \tilde{l}(\mathfrak{D}_{\mathcal{L}_m}).$$

Будем считать, что объем выборки достаточный, если:

$$\forall m_1, m_2 > m^* \quad |l(m_1) - l(m_2)| < \varepsilon,$$

где ε - достаточно малое пороговое значение.

Для оценки вектора параметров используется принцип максимума правдоподобия:

$$\hat{\mathbf{w}} = \arg \max_{\mathbf{w} \in \mathbb{W}} L(\mathfrak{D}_{\mathcal{L}_m}, \mathbf{w}).$$

Рассмотрим матрицу информации Фишера:

$$\mathbf{I}(\mathfrak{D}, \mathbf{w}) = -\nabla \nabla l(\mathfrak{D}, \mathbf{w}).$$

Будем считать, что $\hat{\mathbf{w}}$ имеет следующее распределение:

$$\hat{\mathbf{w}} \sim \mathcal{N}(\mathbf{m}, \mathbf{V}), \quad \mathbf{V} = \mathbf{I}^{-1}(\mathfrak{D}, m).$$

Для линейной регрессии $\mathbb{Y} = \mathbb{R}$, где y представимо в виде:

$$y = \mathbf{x}^\top \mathbf{w} + \varepsilon,$$

где $\varepsilon \sim \mathcal{N}(0, 1)$.

Аппроксимация плотности апостериорной вероятности имеет вид:

$$f(y, \mathbf{x}, \mathbf{w}) = \mathcal{N}(y | \mathbf{x}^\top \mathbf{w}, 1).$$

Для логистической регрессии $\mathbb{Y} = \{0, 1\}$, где y является бернуллиевской случайной величиной:

$$y \sim \mathcal{Be}(\theta),$$

где θ - неизвестный параметр распределения. Аппроксимация плотности апостериорной вероятности имеет вид:

$$f(y, \mathbf{x}, \mathbf{w}) = \mathcal{Be}(y | \theta), \quad \theta = \frac{1}{1 + \exp(-\mathbf{x}^\top \mathbf{w})}.$$

3 Анализ эффективности модели

Методы байесовских оценок объема выборки основаны на ограничении некоторой выбранной характеристики модели. Для анализа эффективности вводится функция от объема выборки, увеличение значений которой интерпретируется как уменьшение эффективности модели. Объем выборки m^* выбирается таким, при котором исследуемая функция не превышает некоторого порогового значения ε .

3.1 Критерий средней длины

Пусть $A(\mathfrak{D}) \subset \mathbb{R}^n$ — множество значений параметров модели \mathbf{w} :

$$A(\mathfrak{D}) = \{\mathbf{w} : \|\mathbf{w} - \hat{\mathbf{w}}\| \leq \alpha\}$$

такое, что

$$P(A(\mathfrak{D})) = 1 - \alpha,$$

где α — некоторое малое значение.

Критерий средней длины выглядит следующим образом:

$$\forall m \geq m^* \quad E_{\mathfrak{D}_m} r_m \leq l,$$

где r_m — радиус шара $A(\mathfrak{D}_m)$, l — некоторое наперед заданное достаточно малое значение.

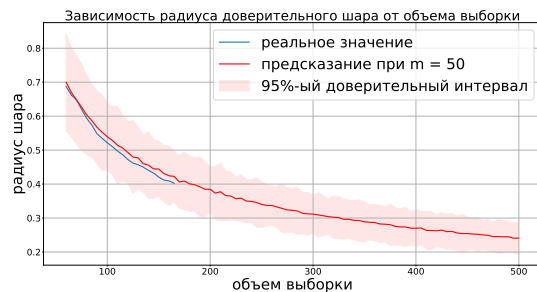
4 Вычислительный эксперимент

Для анализа точности и эффективности предлагаемого метода был проведен вычислительный эксперимент. В качестве данных использовались выборки, описанные в таблице 1.

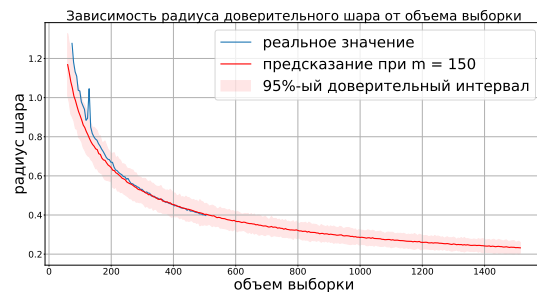
Таблица 1 Описание выборок

Выборка	Тип задачи	Размер выборки	Число признаков
Servo	регрессия	167	4
Boston	регрессия	506	13
Diabetes	регрессия	442	5

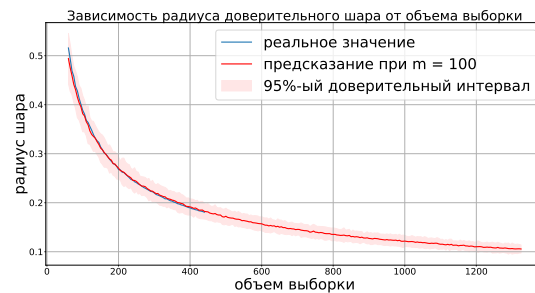
В ходе эксперимента был модифицирован критерий средней длины для линейной регрессии, а именно была построена аппроксимация функции эффективности при большем объеме выборки при помощи аппроксимации ковариационной матрицы вектора параметров через матрицу информации Фишера.



(а) Servo



(б) Boston



(в) Diabetes

Рис. 1 ALC метод

На графике 1 показана зависимость значения функции эффективности от объема выборки при разных данных. Синим цветом обозначено посчитанное значение при данном объеме, красным — аппроксимация при подвыборке фиксированного размера. Реальное значение функции эффективности попадает в доверительный интервал, что говорит о работоспособности метода.

Литература

- [1] *S. G. Self and R. H. Mauritsen* Power/sample size calculations for generalized linear models // Biometrics, 1988.
- [2] *G. Shieh* On power and sample size calculations for likelihood ratio tests in generalized linear models // Biometrics, 2000.

-
- [3] *G. Shieh* On power and sample size calculations for Wald tests in generalized linear models // Journal of Statistical Planning and Inference, 2005.
 - [4] *D. B. Rubin and H. S. Stern* Sample size determination using posterior predictive distributions // Sankhya : The Indian Journal of Statistics Special Issue on Bayesian Analysis, 1998.
 - [5] *Maher Qumsiyeh* Using the bootstrap for estimation the sample size in statistical experiments // Journal of modern applied statistical methods, 2002.