

Численные методы оценки оптимального объёма выборки для логистической и линейной регрессии

Гадаев Тамаз

Московский Физико-Технический Институт
Физтех-школа прикладной математики и информатики
Кафедра Интеллектуальные системы
Научный руководитель: д. ф.-м. н. Стрижов В. В.

Математические методы распознавания образов - 2019

28 ноября 2019 г.

Проблема

Завышенные оценки оптимального объёма выборки приводят к потерям ресурсов в случаях, когда сбор выборки затратен

Цель

Точно оценить оптимальный объём выборки для построения обобщённо-линейной модели

Предлагается

Провести исследование методов на основе статистического анализа мощности, байесовского вывода и эвристик

Методы основанные на анализе мощности статистических критериев

- S. G. Self and R. H. Mauritsen Power/sample size calculations for generalized linear models // Biometrics, 1988. Vol. 44. P. 79–86.
- G. Shieh On power and sample size calculations for likelihood ratio tests in generalized linear models // Biometrics, 2000. Vol. 56. P. 1192–1196

Методы основанные на Байесовском выводе

- L. Joseph and R. du Berger and P. Be'lisle Bayesian and mixed bayesian likelihood criteria for sample size determination // Statistician, 1995. Vol. 16. P. 769–781.
- A. Motrenko and V. Strijov and W. Weber Sample size determination for logistic regression // Journal of Computational and Applied Mathematics, 2014. Vol. 255. P. 743–752.
- 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. Vol. 60. P. 161–175.

Эвристичекие методы

- Maher Qumsiyeh Using the bootstrap for estimation the sample size in statisticalexperiments // Journal of modern applied statistical methods, 2002. Vol. 8(3). P. 305P321.

Дана выборка

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

по \mathfrak{D}_m восстанавливается зависимость $f(\mathbf{x}) = y$, где

$$f(\mathbf{x}) = \begin{cases} \mathbf{x}\mathbf{w}^T, & \text{если } y \in \mathbb{R}, \quad (\text{линейная регрессия}) \\ \mathbb{I}[\sigma(\mathbf{x}\mathbf{w}^T) > 0.5], & \text{если } y \in \{0, 1\}, \quad (\text{логистическая регрессия}) \end{cases}$$

Оценка вектора параметров $\hat{\mathbf{w}} = \arg \max_{\mathbf{w}} \ln L(\mathfrak{D}, \mathbf{w})$, где

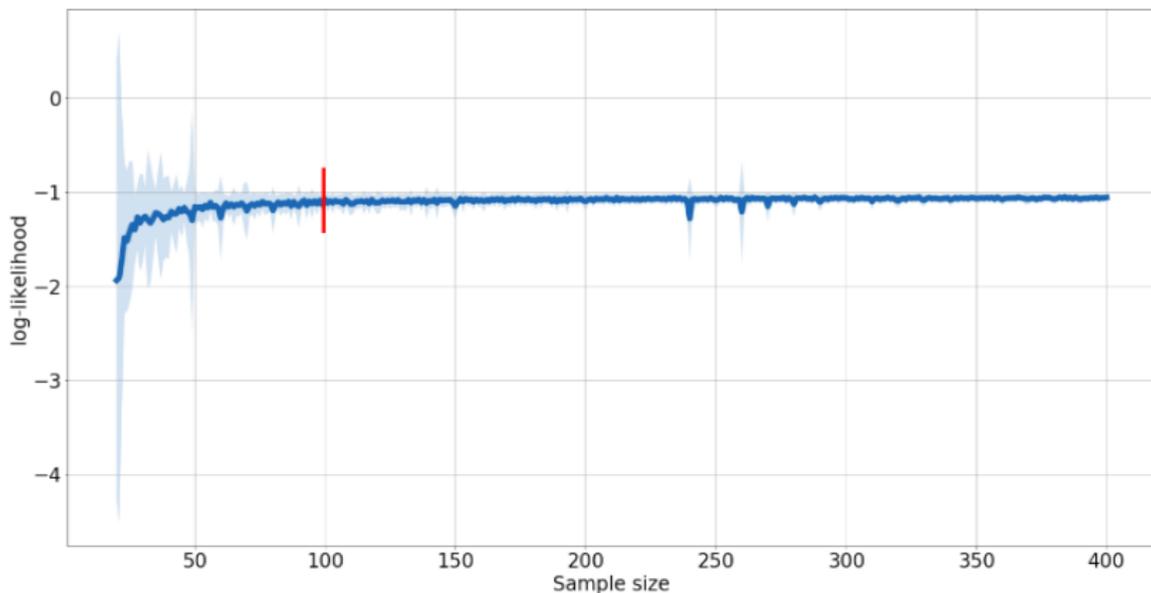
$$\ln L(\mathfrak{D}, \mathbf{w}) = \begin{cases} \sum_{i=1}^m (y_i - f(\mathbf{x}_i))^2, & \text{если } y \in \mathbb{R} \\ \sum_{i=1}^m y_i \ln(f(\mathbf{x}_i)) + (1 - y_i) \ln(1 - f(\mathbf{x}_i)), & \text{если } y \in \{0, 1\} \end{cases}$$

Задача

Необходимо адекватно определить оптимальность объёма выборки для нахождения $\hat{\mathbf{w}}$ и в соответствии с этим определением сделать точную оценку m^*

Пример: оценка по дисперсии логарифма правдоподобия

В этом примере оптимальным выбирается такой размер выборки, при котором дисперсия логарифма правдоподобия выходит на константу.



Модель порождения данных:

$$p(y|\mathbf{u}, \mathbf{v}, \mathbf{w}_u, \mathbf{w}_v) = \exp\left(\frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi)\right)$$

Вектор признаков \mathbf{x} представлен как $\mathbf{x} = [\mathbf{u}, \mathbf{v}]$, вектор параметров $\mathbf{w} = [\mathbf{w}_u, \mathbf{w}_v]$, истинный вектор параметров $\mathbf{m} = [\mathbf{m}_u, \mathbf{m}_v]$.

Рассмотрим гипотезу:

$$H_0 : \mathbf{m}_u = \mathbf{m}_u^0, \quad H_1 : \mathbf{m}_u \neq \mathbf{m}_u^0$$

где \mathbf{m}_u^0 — гиперпараметр метода.

Основная идея

В каждом из статистических методов вычисляется статистика $T(\mathcal{D}_m, \mathbf{w}_u^0)$. Известно, что $T \sim \chi_k^2$ при H_0 и $T \sim \chi_k^2(\gamma)$ при H_1 , где γ в каждом методе вычисляется по-своему.

Условие оптимальности:

$$[P(H_0 \text{ отвергнута} | H_0) = \alpha] \cap [P(H_0 \text{ отвергнута} | H_1) = \beta]$$

Найти m^* можно решив уравнение

$$\chi_{k,1-\alpha}^2 = \chi_{k,\beta}^2(\gamma(m))$$

Пусть статистики $S_{m,u}(\mathbf{w}_u, \mathbf{w}_v)$ и $S_{m,v}(\mathbf{w}_u, \mathbf{w}_v)$ - производные логарифма правдоподобия выборки \mathfrak{D}_m по параметрам \mathbf{w}_u и \mathbf{w}_v соответственно. Рассмотрим $\mathbf{s}_m = S_{m,u}(\mathbf{m}_u^0, \hat{\mathbf{w}}_v^0)$, где $\hat{\mathbf{w}}_v^0$ определяется из уравнения

$$S_{m,v}(\mathbf{m}_u^0, \mathbf{w}_v) = 0$$

Метод множителей Лагранжа использует статистику $LM = \mathbf{s}_m^T \mathbf{Q}_m^{-1} \mathbf{s}_m$, где \mathbf{Q}_m — матрица ковариации \mathbf{s}_m .

$LM \stackrel{H_0}{\sim} \chi_k^2$, где k размер вектора \mathbf{m}_u

$LM \stackrel{H_1}{\sim} \chi_k^2(\gamma)$ [1], где γ — параметр нецентральности, задаваемый в следующем виде:

$$\gamma = \boldsymbol{\xi}_m^T \boldsymbol{\Sigma}_m^{-1} \boldsymbol{\xi}_m = m \boldsymbol{\xi}^T \boldsymbol{\Sigma}^{-1} \boldsymbol{\xi} = m \gamma^0,$$

где $\boldsymbol{\xi}_m = \mathbb{E} \mathbf{s}_m$, $\boldsymbol{\Sigma}_m = \mathbb{E}[\mathbf{s}'_m \mathbf{s}_m^T]$.

Также,

$$\gamma^* : \chi_{k,1-\alpha}^2 = \chi_{k,\beta}^2(\gamma)$$

Тогда

$$m^* = \frac{\gamma^*}{\gamma^0}$$

Рассмотрим статистику логарифма отношения правдоподобий:

$$LR = 2 \left(l(\mathcal{D}, \hat{\mathbf{w}}) - l(\mathcal{D}, \hat{\mathbf{w}}^0) \right)$$

где $\hat{\mathbf{w}} = [\hat{\mathbf{w}}_u, \hat{\mathbf{w}}_v]$ — оценка максимума правдоподобия, $\hat{\mathbf{w}}^0 = [\mathbf{m}_u^0, \hat{\mathbf{w}}_v^0]$ — вектор параметров модели, который максимизирует правдоподобие при фиксированном \mathbf{m}_u^0 .

$LR \stackrel{H_0}{\sim} \chi_k^2$, где k размер вектора \mathbf{m}_u

$LR \stackrel{H_1}{\sim} \chi_k^2 [2]$, где γ — параметр нецентральности:

$$\gamma = m\Delta^*, \quad \Delta^* = \mathbb{E} \left[2a^{-1}(\phi) \{(\theta - \theta^*) \nabla b(\theta) - b(\theta) + b(\theta^*)\} \right],$$

где параметры θ и θ^* считаются по параметрам $\mathbf{w} = [\mathbf{w}_u, \mathbf{w}_v]$ и $\mathbf{w}^* = [\mathbf{w}_u^0, \mathbf{w}_v^*]$ соответственно, вектор параметров \mathbf{w}_v^* задается как решение следующего уравнения:

$$\lim_{m \rightarrow \infty} m^{-1} \mathbb{E} \left[\frac{\partial l(\mathcal{D}, [\mathbf{m}_u^0, \mathbf{w}_v])}{\partial \mathbf{w}_v} \right] = 0.$$

Тогда при заданных α и β получаем оценку размера выборки m^* :

$$m^* = \frac{\gamma^*}{\Delta^*}, \quad \gamma^* : \chi_{k,1-\alpha}^2 = \chi_{k,\beta}^2(\gamma)$$

Вместо обычной логистической и линейной регрессии рассмотрим байесовскую.

Апостериорное распределение:

$$p(w|\mathcal{D}) = \mathcal{N}(w|\hat{w}, V)$$

где $V = I^{-1}(\mathcal{D}, \hat{w}) = -(\nabla\nabla L(\mathcal{D}, \hat{w}))^{-1}$ — матрица, обратная матрице Фишера

Для линейной регрессии $\hat{w} = (\mathbf{X}^T \mathbf{X} + \alpha \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y}$

Для логистической регрессии воспользуемся аппроксимацией Лапласа в точке МП $\hat{w} = \mathit{argmax}_w L(\mathcal{D}, w)$

Основная идея

Вычисляется апостериорное распределение для разных подвыборок с возрастающим размером до тех пор, пока не будет выполнено некоторое условие

Условие оптимальности: Апостериорное распределение достаточно локализовано. Локализация может быть описана тремя способами (среднее покрытие, средняя длина, средняя дисперсия).

Критерий средней апостериорной дисперсии

Пусть $D\hat{\mathbf{w}}$ — Дисперсия апостериорного распределения вектора параметров $\hat{\mathbf{w}}$. Тогда $E_{\mathcal{D}_m} D[\hat{\mathbf{w}}|\mathcal{D}_m]$ — матожидание этой дисперсии по всем выборкам размера m из одного распределения

Размер выборки m^* находится из условия:

$$\forall m \geq m^* E_{\mathcal{D}_m} D[\hat{\mathbf{w}}|\mathcal{D}_m] \leq l$$

Критерий среднего покрытия

Пусть $A(\mathcal{D}) \subset \mathbb{R}^n$ некоторое множество значений параметров модели \mathbf{w} :
 $A(\mathcal{D}) = \{\mathbf{w} : \|\mathbf{w} - \hat{\mathbf{w}}\| \leq l\}$, где l — заданный радиус шара.

Критерий для определения m^* :

$$\forall m \geq m^* E_{\mathcal{D}_m} P\{\mathbf{w} \in A(\mathcal{D}_m)\} \geq 1 - \alpha$$

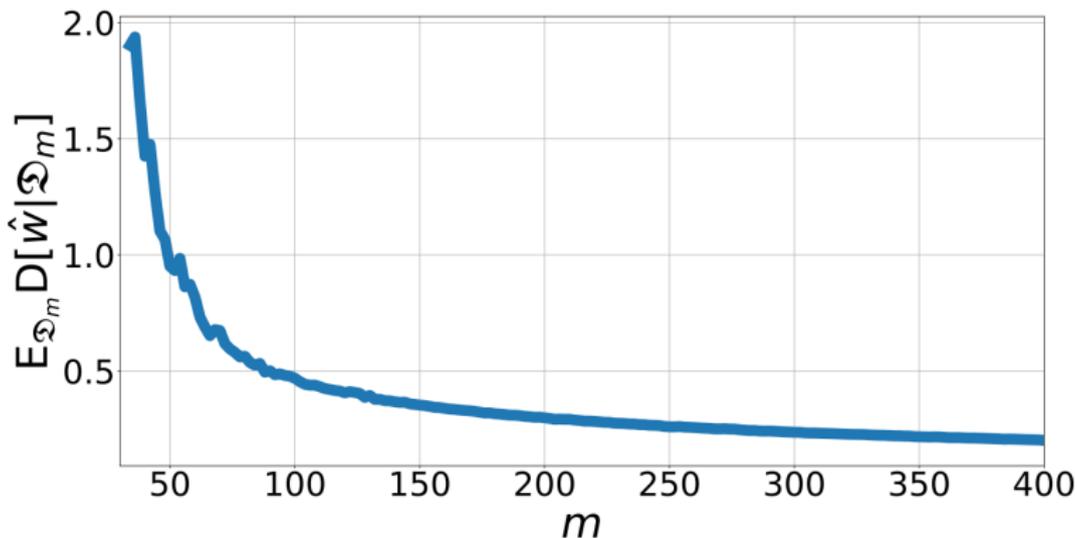
Критерий средней длины

Аналогично предыдущему, но $A(\mathcal{D})$ определяется как $P(A(\mathcal{D})) = 1 - \alpha$

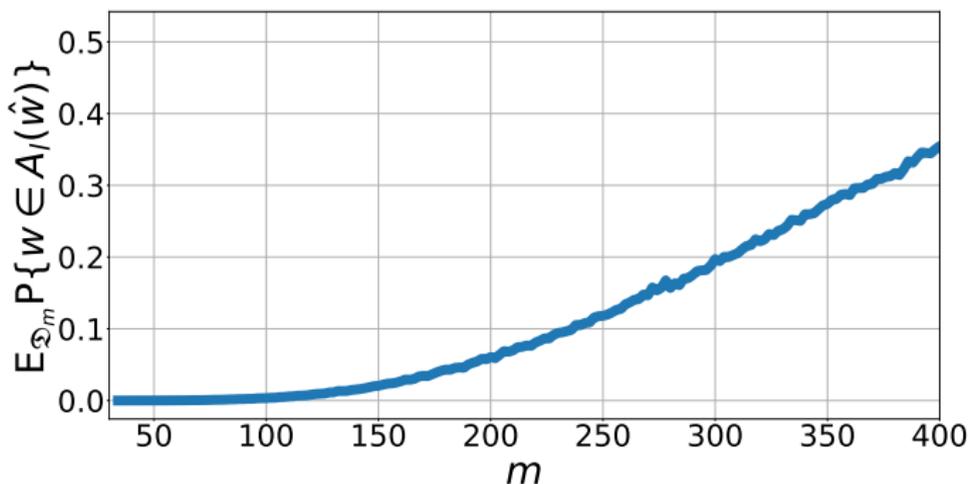
Пусть r_m радиус шара $A(\mathcal{D}_m)$. Критерий для определения m^* :

$$\forall m \geq m^* E_{\mathcal{D}_m} r_m \leq l$$

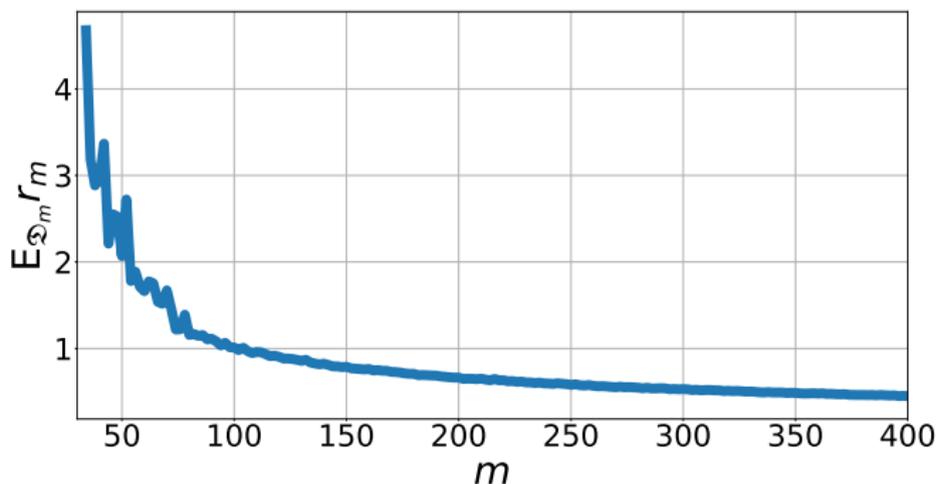
Локализация апостериорного распределения параметров $p(\hat{\mathbf{w}}|\mathcal{D}_m)$ как уменьшение средней дисперсии $E_{\mathcal{D}_m} D[\hat{\mathbf{w}}|\mathcal{D}_m]$ по бутстрепным выборкам.



Локализация апостериорного распределения параметров $p(\hat{w}|\mathcal{D}_m)$ как увеличение средней вероятности попадания вектора параметров в шар радиуса l . Усреднение производится по бутстрепным выборкам.



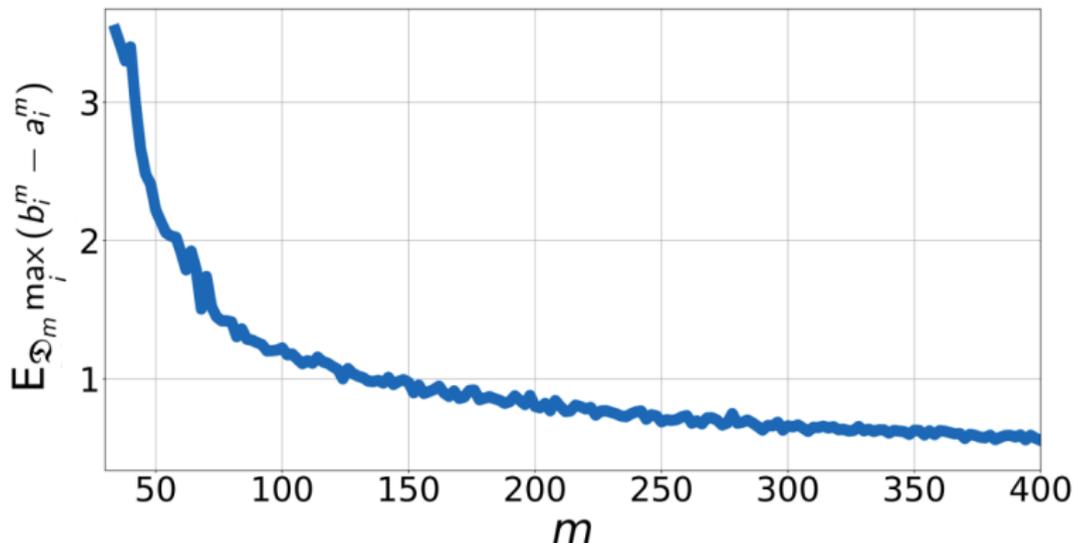
Локализация апостериорного распределения параметров $p(\hat{w}|\mathcal{D}_m)$ как среднего уменьшение радиуса 0.95-доверительного шара по бутстрепным выборкам.



По выборке размера m построим квантильные бутстрепные доверительные интервалы $(a_1^m, b_1^m), (a_2^m, b_2^m), \dots, (a_n^m, b_n^m)$ уровня значимости α для каждого параметра модели.

Оптимальный объём выборки m^* находится из условия малой длины этих интервалов:

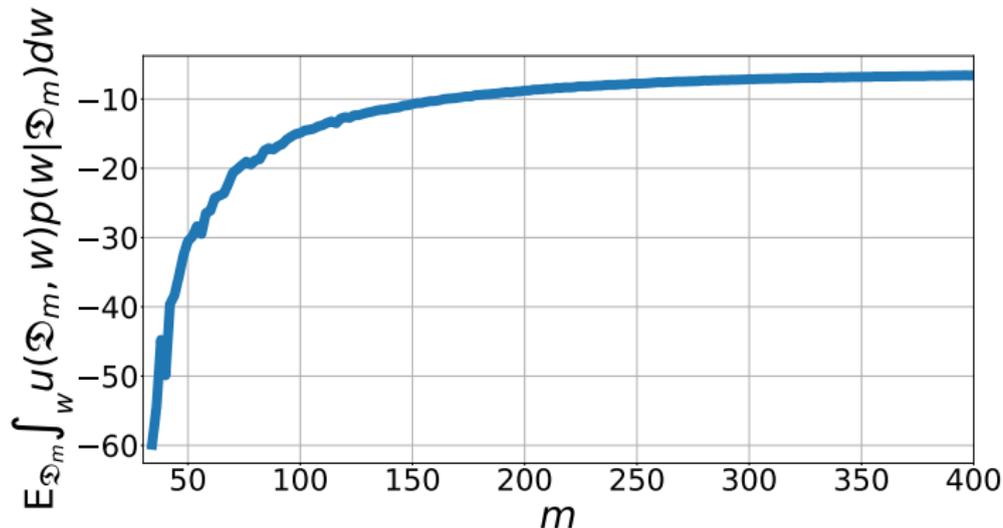
$$m^* : \forall m \geq m^* \max_i (b_i^m - a_i^m) < l.$$



Введём функцию полезности модели $u(\mathcal{D}_m)$ (в данном случае $u(\mathcal{D}_m) = \ln L(\mathcal{D}_m, \hat{w})$) и будем искать

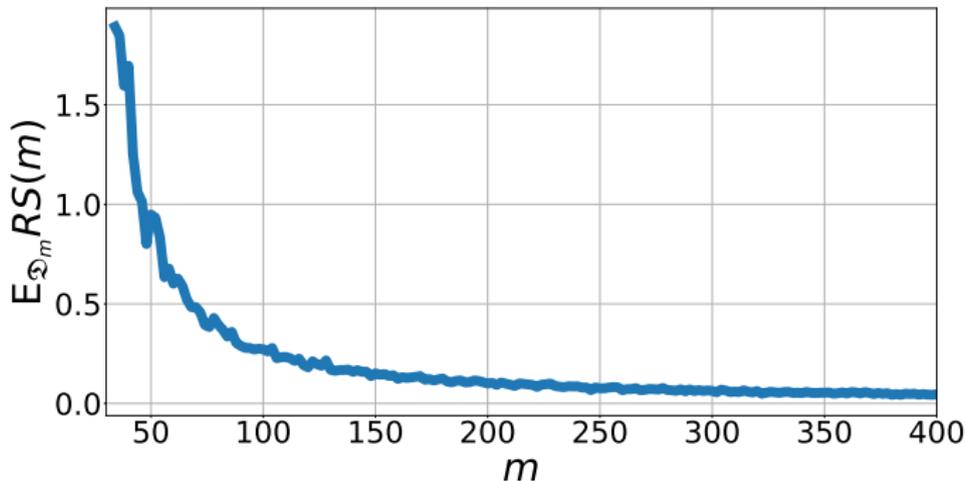
$$m^* = \arg \max_m (\mathbb{E}_{\mathcal{D}_m} [u(\mathcal{D}_m)] - cm),$$

где c - коэффициент штрафа за размер выборки

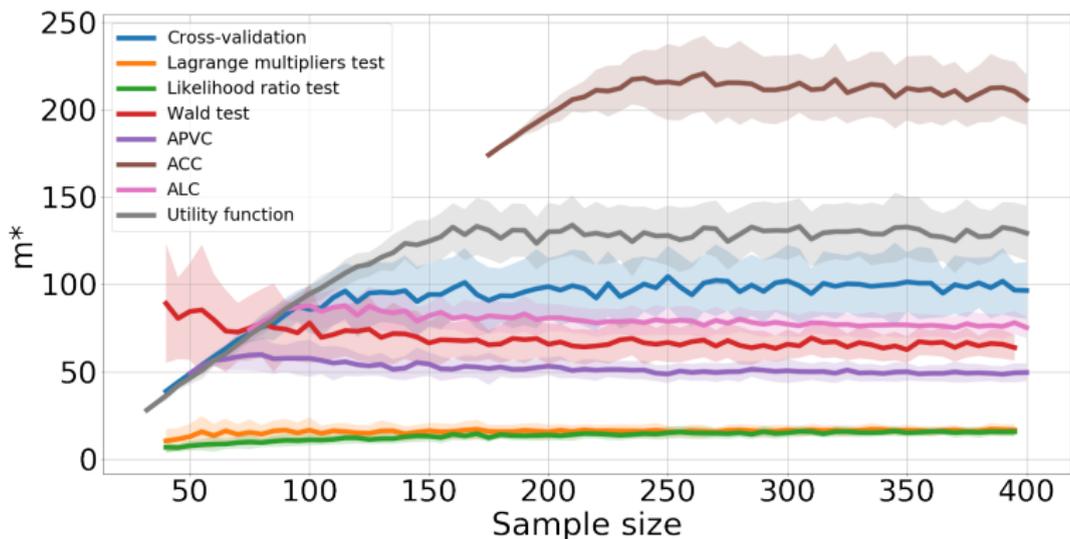


Размер выборки m^* определяется так что для любого $m \geq m^*$ выполнено $RS(m) \geq 1 - \varepsilon$, где ε выбирается малым.

При этом $RS(m) = \ln \frac{L(\mathcal{D}_{\mathcal{L}(m)}, \hat{\mathbf{w}})}{L(\mathcal{D}_{\mathcal{T}(m)}, \hat{\mathbf{w}})}$, где \mathcal{T}, \mathcal{L} — обучающая и тестовая выборки для подвыборки размера m .



Дисперсия при каждом m вычисляется по разным подвыборкам размера m . Все графики выходят на константу.



Выводы

- Исследовано поведение различных методов оценки оптимального объёма выборки
- Описана методология оценки оптимального объёма выборки
- Реализовано программное обеспечение для ретроспективной и предсказательной оценки оптимального объёма выборки

-  *S. G. Self and R. H. Mauritsen* Power/sample size calculations for generalized linear models // *Biometrics*, 1988. Vol. 44. P. 79–86.
-  *G. Shieh* On power and sample size calculations for likelihood ratio tests in generalized linear models // *Biometrics*, 2000. Vol. 56. P. 1192–1196.
-  *G. Shieh* On power and sample size calculations for Wald tests in generalized linear models // *Journal of Statistical Planning and Inference*, 2005. Vol. 128. P. 43–59.
-  *E. Demidenko* Sample size determination for logistic regression revisited // *Statist. Med.*, 2007. Vol. 26. P. 3385–3397.
-  *A. Motrenko and V. Strijov and W. Weber* Sample size determination for logistic regression // *Journal of Computational and Applied Mathematics*, 2014. Vol. 255. P. 743–752.
-  *Maher Qumsiyeh* Using the bootstrap for estimation the sample size in statistical experiments // *Journal of modern applied statistical methods*, 2002. Vol. 8(3). P. 305–321.
-  *Fei Wang and Alan E. Gelfand* A Simulation-based Approach to Bayesian Sample Size Determination for Performance under a Given Model and for Separating Models // *Statistical Science*, 2002. Vol. 17. P. 193–208.
-  *D. V. Lindley* The choice of sample size // *The Statistician*, 1997. Vol. 46. P. 129–138.