

Федеральное государственное автономное образовательное учреждение
высшего образования
«Московский физико-технический институт (национальный исследовательский
университет)»

Направление подготовки:

Направленность (профиль) подготовки:

Студент:

Бучнев Валентин Сергеевич

(подпись студента)

Научный руководитель:

Стрижов Вадим Викторович

(подпись научного руководителя)

Москва 2020

Содержание

1. Введение	4
2. Постановка задачи	6
2.1. Вычисление функции $l(m)$	8
2.2. Аппроксимация $\phi(m) \sim \hat{l}(m)$	8
3. Вычислительный эксперимент	10
3.1. Эксперимент на синтетических выборках	10
3.2. Эксперимент на выборках из UCI репозитория	15
4. Заключение	19

Аннотация

Исследована проблема снижения затрат на сбор данных, необходимых для построения адекватной модели. Рассматриваются задачи линейной и логистической регрессий. Для решения этих задач требуется, чтобы выборка содержала достаточное 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, бутстреп.*

1. Введение

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

Существует большое количество методов оценки достаточного объёма выборки. Например, тест множителей Лагранжа, тест отношения правдоподобия и тест Вальда. В работах [1, 2, 3] на основе данных методов построена оценка достаточного объёма выборки. Основным минус этих методов заключается в том, что статистики, используемые в критериях, имеют асимптотическое распределение и требуют большого объёма выборки.

Существуют байесовские оценки объёма выборки: критерий средней апостериорной дисперсии, критерий среднего покрытия, критерий средней длины и метод максимизации полезности. Первые три метода, описанные в работе [4], требуют анализа некоторой функции эффективности от размера выборки. Используя некоторое решающее правило, по данной функции определяется достаточный объём выборки. Главный минус этих методов заключается в том, что они не позволяют построить аппроксимацию функции эффективности при большем объёме данных. Метод максимизации полезности максимизирует ожидание некоторой функции полезности по объёму выборки. Все эти методы опираются на апостериорное распределение, что требует достаточно большого объёма выборки.

Существуют также модели, аппроксимирующие зависимость функции ошибки от объёма выборки и предсказывающие достаточный объём выборки. В работах [5, 6] описаны такие модели для решения определённой задачи классификации на конкретном наборе данных. В работе [6] из датасета извлекаются характерные для задачи признаки (дисперсия, количество различных генов и т.д.), которые далее используются для построения аппроксимирующей модели.

Предлагается исследовать зависимость значения логарифма правдоподобия от размера доступной выборки, а также его дисперсию. В данной работе используется модель, предсказывающая значение логарифма правдоподобия при большем объёме выборки. Для вычислительного эксперимента предлага-

ется использовать классические выборки из UCI репозитория и синтетические данные.

2. Постановка задачи

Дана выборка размера m :

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

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

Выборка является простой: элементы порождены независимо из одного распределения с фиксированными неизвестными параметрами, вероятность попадания каждого элемента в выборку одинакова. Предполагается, что выборка \mathfrak{D} порождена согласно следующей гипотезе: модель, порождающая данные, задается в следующем виде:

$$y_i = f(\mathbf{x}_i, \mathbf{w}, \beta), \quad (1)$$

где $\mathbf{w} \in \mathbb{W}$ — вектор параметров, β — дисперсия зависимой переменной. Зависимая переменная y аппроксимируется обобщенно линейной моделью:

$$\hat{y}_i = f(\mathbf{x}_i, \mathbf{w}) = \mu(\mathbf{w}^\top \mathbf{x}_i), \quad (2)$$

где μ — функция связи, для модели линейной регрессии:

$$\mu = id, \quad (3)$$

для логистической регрессии:

$$\mu(\mathbf{w}^\top \mathbf{x}_i) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_i)}. \quad (4)$$

Предполагается, что при восстановлении параметров линейной регрессии (3), зависимая переменная порождается нормальным распределением:

$$p(y|\mathbf{x}, \mathbf{w}) \sim \mathcal{N}(f(\mathbf{x}, \mathbf{w}), \hat{\beta}),$$

где $\hat{\beta}$ — выборочная дисперсия зависимой переменной y . Для модели логистической регрессии зависимая переменная порождается бернуллиевским распределением:

$$p(y|\mathbf{x}, \mathbf{w}) \sim \mathcal{Be}(f(\mathbf{x}, \mathbf{w})).$$

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

$$L(\mathbf{w}|\mathfrak{D}) = \prod_{y, \mathbf{x} \in \mathfrak{D}} p(y|\mathbf{x}, \mathbf{w}), \quad l(\mathbf{w}|\mathfrak{D}) = \sum_{y, \mathbf{x} \in \mathfrak{D}} \log p(y|\mathbf{x}, \mathbf{w}). \quad (5)$$

Для множества индексов \mathcal{A} независимой переменной $\mathbf{x} = [x_1, \dots, x_n]^\top$, в векторы $\mathbf{x}_{\mathcal{A}}, \mathbf{w}_{\mathcal{A}}$ входят только те элементы, индексы которых принадлежат \mathcal{A} . Определим выборку $\mathfrak{D}_{\mathcal{A}}$ как множество независимых переменных $\mathbf{x}_{\mathcal{A}}$ и соответствующих им зависимых переменных y .

Для получения оптимального набора параметров и оценки вектора параметров используется принцип максимума правдоподобия:

$$\hat{\mathbf{w}}, \hat{\mathcal{A}} = \arg \max_{\mathbf{w} \in \mathbb{W}, \mathcal{A} \subseteq \mathbb{J}} L(\mathbf{w}_{\mathcal{A}}|\mathfrak{D}_{\mathcal{A}}), \quad (6)$$

где $\mathbb{J} = \{1, 2, \dots, n\}$ — множество индексов.

В качестве функции ошибки используются следующие функции:

$$\begin{aligned} S_{\text{reg}}(\mathbf{w}|\mathfrak{D}) &= \frac{1}{|\mathfrak{D}|} \sum_{\mathbf{x}, y \in \mathfrak{D}} (y - f(\mathbf{x}, \mathbf{w}))^2, \\ S_{\text{class}}(\mathbf{w}|\mathfrak{D}) &= \frac{1}{|\mathfrak{D}|} \sum_{\mathbf{x}, y \in \mathfrak{D}} (y \ln f(\mathbf{x}, \mathbf{w}) + (1 - y) \ln(1 - f(\mathbf{x}, \mathbf{w}))). \end{aligned} \quad (7)$$

Определим функцию обратной экспоненты от функции ошибки (7):

$$e^{-S(\mathbf{w}|\mathfrak{D})}. \quad (8)$$

Заметим, что задача (6) эквивалентна задаче:

$$\hat{\mathbf{w}}, \hat{\mathcal{A}} = \arg \max_{\mathbf{w} \in \mathbb{W}, \mathcal{A} \subseteq \mathbb{J}} e^{-S(\mathbf{w}|\mathfrak{D}_{\mathcal{A}})}, \quad (9)$$

где S — функция ошибки (7).

Разделим выборку \mathfrak{D} на обучающую и тестовую:

$$\mathfrak{D} = \mathfrak{D}_{\mathcal{L}} \sqcup \mathfrak{D}_{\mathcal{T}} \quad (10)$$

Требуется по начально заданной выборке размера $m_0 \ll m$ получить оценку минимального достаточного объёма m^* . Для введения понятия достаточности объёма выборки рассмотрим ожидаемое значение функции (8) по разным

подвыборкам $\mathfrak{D}'_{\mathcal{L}}, \mathfrak{D}'_{\mathcal{T}}$ размера m' обучающей и тестовой выборки (10). Оценка вектора параметров $\hat{\mathbf{w}}$ является решением задачи (9) для выборки $\mathfrak{D}'_{\mathcal{L}}$:

$$l(m') = m'^{-1} \sum_{\substack{\mathfrak{D}'_{\mathcal{L}} \subset \mathfrak{D}_{\mathcal{L}} \\ \mathfrak{D}'_{\mathcal{T}} \subset \mathfrak{D}_{\mathcal{T}} \\ |\mathfrak{D}'_{\mathcal{L}}| = |\mathfrak{D}'_{\mathcal{T}}| = m'}} e^{-S(\hat{\mathbf{w}}(\mathfrak{D}'_{\mathcal{L}}) | \mathfrak{D}'_{\mathcal{T}})}. \quad (11)$$

Будем считать, что объем выборки m^* достаточен, если:

$$\forall m' > m^* \quad l(m') > (1 - \delta) \max_{m > m^*} l(m),$$

где δ — достаточно малое пороговое значение.

2.1. Вычисление функции $l(m)$

Для получения приближённого значения функции $l(m')$ введем процедуру бутстрепа:

- 1) Равновероятно генерируются случайные подвыборки $\mathfrak{D}'_{\mathcal{L}}, \mathfrak{D}'_{\mathcal{T}}$ размера m' :
 - $\mathfrak{D}' \sim U(\mathfrak{D})$ — вариант с полной информацией,
 - $\mathfrak{D}' \sim U(\mathfrak{D}^0), \mathfrak{D}^0 \subset \mathfrak{D}, |\mathfrak{D}^0| = m'$ — вариант с неполной информацией.
- 2) Для полученных подвыборок вычисляется значение функции (8):
- 3) пп. 1-2 повторяются K раз, оценка $\hat{l}(m')$ функции $l(m')$ равняется среднему арифметическому среди всех полученных значений функции (8) на всех итерациях:

2.2. Аппроксимация $\phi(m) \sim \hat{l}(m)$

Для предсказания значения функции $\hat{l}(m)$ при $m > m_0$ введем параметрическое семейство функций

$$\Phi = \{\phi(m) = a + b \cdot e^{c \cdot m} \mid a, b \in \mathbb{R}, c \in (-\infty, 0)\} \quad (12)$$

Аппроксимация функции $\hat{l}(m)$ является решением следующей задачи:

$$\hat{\phi} = \arg \min_{\phi \in \Phi} \text{MAE}(\hat{l}, \phi, 1, m_0), \quad (13)$$

где MAE — средняя абсолютная ошибка:

$$\text{MAE}(\psi, \phi, m_1, m_2) = \frac{1}{m_2 - m_1 + 1} \sum_{i=m_1}^{m_2} |\phi(i) - \psi(i)|. \quad (14)$$

3. Вычислительный эксперимент

Для отбора нужного количества признаков используем Лассо регрессию с функцией ошибки:

$$MSE(\mathbf{x}, y, \mathbf{w}) = \frac{1}{m} \sum_{i=1}^m (y_i - \mathbf{w}^\top \mathbf{x}_i)^2 + \tau \|\mathbf{w}\|_1. \quad (15)$$

Используется тот факт, что Лассо регрессия обнуляет малозначимые признаки, и чем больше коэффициент τ , тем больше признаков обнуляется. Таким образом можно упорядочить признаки по значимости. Для того, чтобы выбрать набор признаков $\mathcal{A}_{n'}$ мощности $|\mathcal{A}_{n'}| = n'$ с наилучшей функцией ошибки, достаточно взять n' самых значимых признаков. Для получения оценки $\hat{\mathbf{w}}$ на наборе признаков $\mathcal{A}_{n'}$ решается оптимизационная задача (6).

3.1. Эксперимент на синтетических выборках

Пусть $\mathbf{X} = [\chi_1, \dots, \chi_n]$ — набор векторов-столбцов. Построим зависимость $\hat{l}(m)$ для различных конфигураций выборок:

- 1) Случайная выборка ($m = 200, n = 10$):

$$\chi_i \sim U(\mathbf{5}, \mathbf{6}), \quad \mathbf{w} \sim U(\mathbf{0}, \mathbf{1}), \quad y = \mathbf{w}^\top \mathbf{X} + \varepsilon,$$

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

- 2) Скоррелированная выборка ($m = 200, n = 10$):

$$\chi_1, \dots, \chi_3 \sim \mathcal{N}(\mathbf{1}, \mathbf{1}), \quad \chi_4, \dots, \chi_6 \sim \chi_0 + \varepsilon, \quad \chi_7, \dots, \chi_{10} \sim \chi_1 + \varepsilon,$$

$$y = 0.3\chi_1 + 0.7\chi_2 + \varepsilon,$$

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

- 3) Ортогональная выборка ($m = 200, n = 10$):

$$\{\chi_i \in \mathbb{R}^m\}_{i=1}^n \text{ — ортогональный набор, } \mathbf{w} \sim U(\mathbf{1}, \mathbf{1}), \quad y = \mathbf{w}^\top \mathbf{X} + \varepsilon,$$

где $\varepsilon \sim \mathcal{N}(0, 0.5)$.

4) Избыточная выборка ($m=200$, $n=10$):

$$\chi_i \sim \mathcal{N}(1, 1),$$

$$w_i \sim U(1, 1), \quad i \leq 5,$$

$$w_i = 0, \quad i > 5,$$

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

где $\varepsilon \sim \mathcal{N}(0, 0.5)$

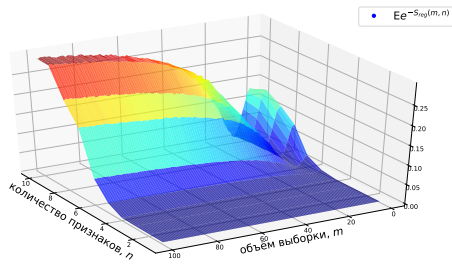
На рис. 1 представлена функция $\hat{l}(m)$, посчитанная с помощью бутстрепа в варианте с полной информацией для различного числа признаков n' . Дисперсия функции $\hat{l}(m)$ монотонно уменьшается с увеличением размера обучающей выборки. Видно, что для всех выборок функция $\hat{l}(m)$ хорошо аппроксимируется семейством функций (12).

В табл. 1 приведены оптимальные значения m^*, n^* для разных конфигураций выборок. Полученные значения n^* ожидаемы из способа генерации. Значения m^* предстоит предсказать по построенной аппроксимации функции $\hat{l}(m)$ при заданном m_0 .

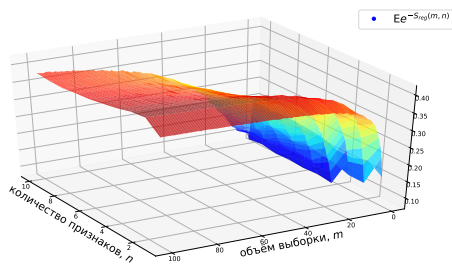
Таблица 1: Оптимальные значения m^*, n^* для различных синтетических выборок

Выборка	m^*	n^*
Случайная выборка	72	10
Скоррелированная выборка	31	2
Ортогональная выборка	45	10
Избыточная выборка	22	5

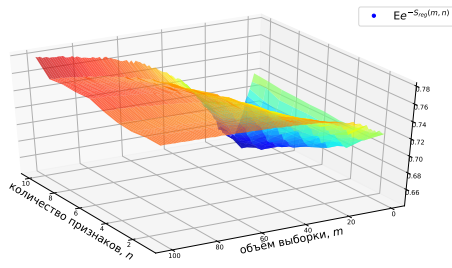
Матожидание



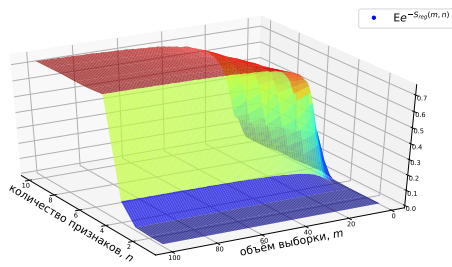
(a) Случайная выборка



(c) Скоррелированная выборка

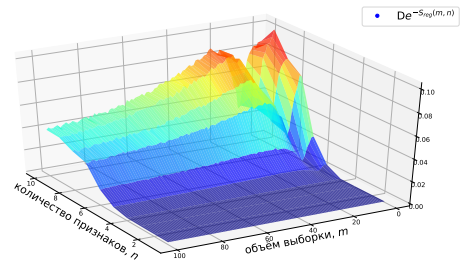


(e) Ортогональная выборка

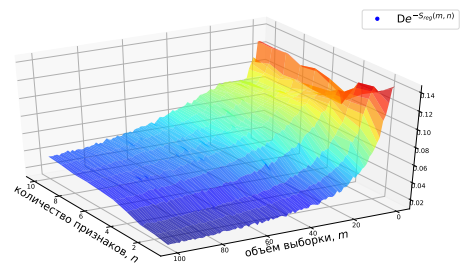


(g) Избыточная выборка

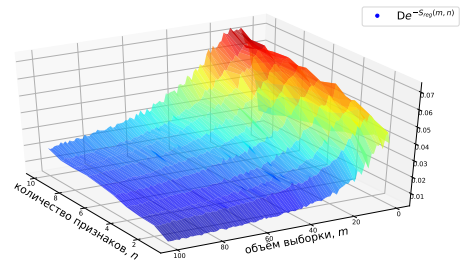
Дисперсия



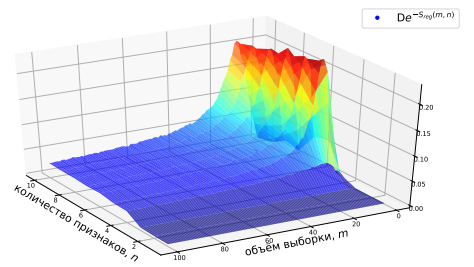
(b) Случайная выборка



(d) Скоррелированная выборка



(f) Ортогональная выборка



(h) Избыточная выборка

Рис. 1: Зависимость значения функции $e^{-S(m,n)}$ от объема выборки m и количества параметров n для синтетических выборок

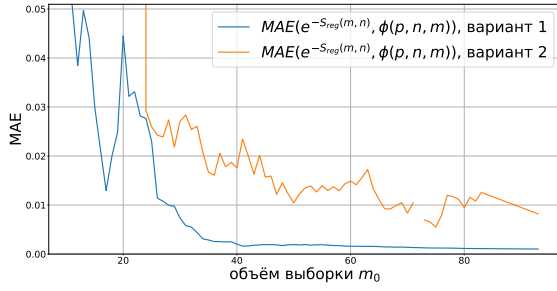
На рис. 2. представлены графики качества аппроксимации функции $\hat{l}(m)$, а также предсказание \hat{m}^* при различных m_0 для разных конфигураций выборок.

Для аппроксимации функции $\hat{l}(m)$ вариант с неполной информацией дает большую ошибку, чем вариант с полной информацией, и это ожидаемый результат. Чем меньше информации, тем хуже качество предсказания.

При $m_0 \rightarrow m$ аппроксимация в варианте с неполной информацией стремится к аппроксимации в варианте с полной информацией, поэтому при больших m_0 предсказания \hat{m}^* для этих двух вариантов практически совпадают.

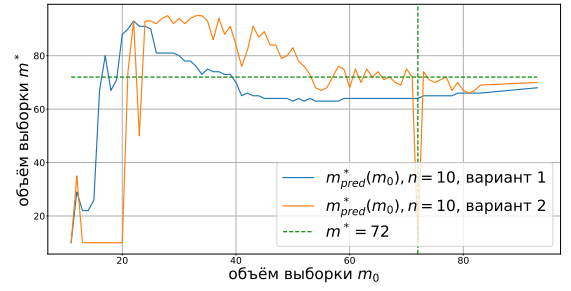
Для случайной и избыточной выборки получилось построить адекватное предсказание при $m_0 < m^*$.

Аппроксимация $e^{-S(m,n)}$

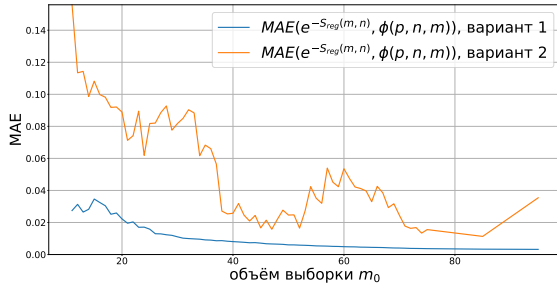


(a) Случайная выборка

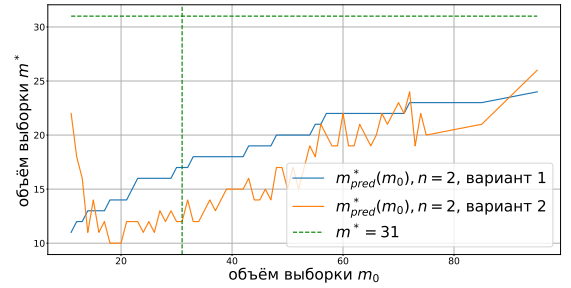
Предсказание m^*



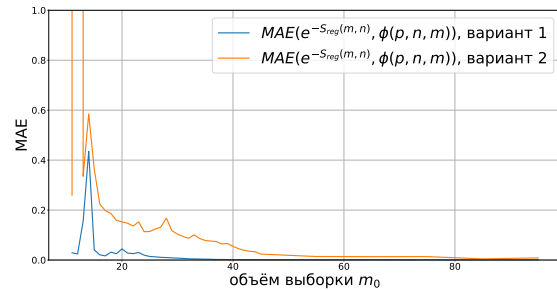
(b) Случайная выборка



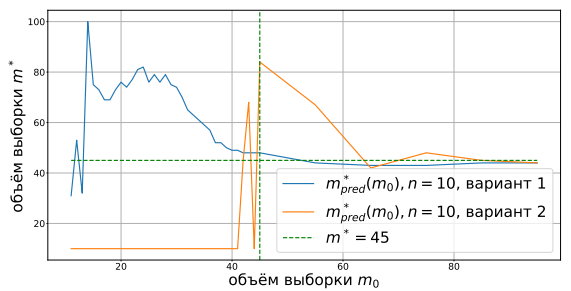
(c) Скоррелированная выборка



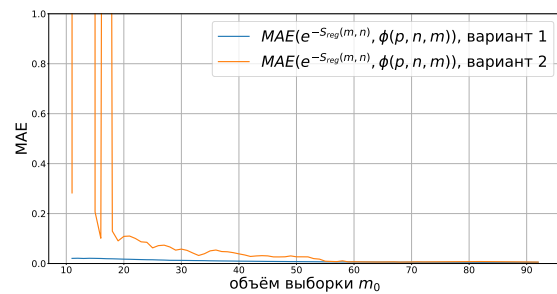
(d) Скоррелированная выборка



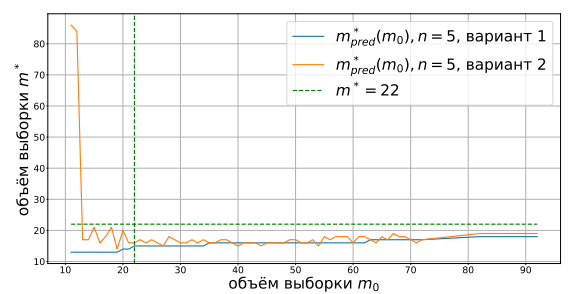
(e) Ортогональная выборка



(f) Ортогональная выборка



(g) Избыточная выборка



(h) Избыточная выборка

Рис. 2: Качество предсказания $e^{-S(m,n)}$ и m^* в зависимости от объема обучающей выборки m_0 для синтетических выборок

3.2. Эксперимент на выборках из UCI репозитория

Для эксперимента использовались выборки из UCI репозитория, описанные в табл. 2.

Таблица 2: Выборки из UCI репозитория

Выборка	Тип задачи	m	n
Diabetes	регрессия	442	11
Boston	регрессия	506	14
Wine	классификация	130	14
Nba	классификация	400	20

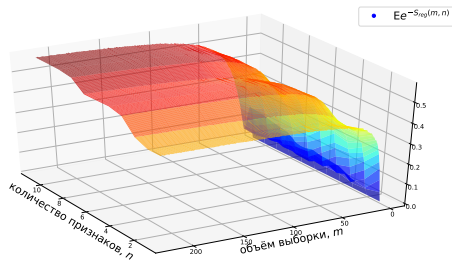
На рис. 3 представлена функция $\hat{l}(m)$, посчитанная с помощью бутстрепа в варианте с полной информацией для различного числа признаков n' . Качественное поведение функции $\hat{l}(m)$ похоже на поведение данной функции для синтетических выборок и также хорошо аппроксимируется семейством функций (12).

В табл. 3 приведены оптимальные значения m^*, n^* для различных выборок. Примечательно, что для задачи классификации достаточный объём выборки меньше, чем для задачи регрессии.

Таблица 3: Оптимальные значения m^*, n^* для различных синтетических выборок

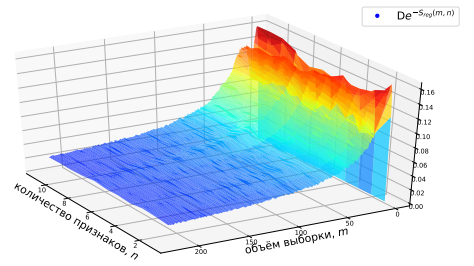
Выборка	m^*	n^*
Diabetes	96	11
Boston	102	14
Wine	27	14
Nba	38	2

Матожидание

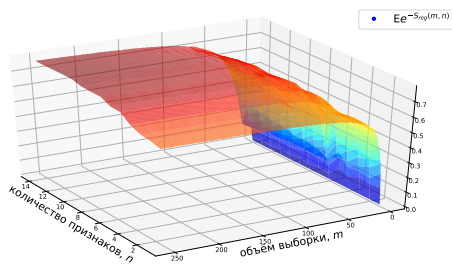


(a) Diabetes

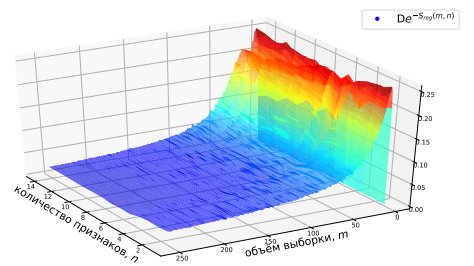
Дисперсия



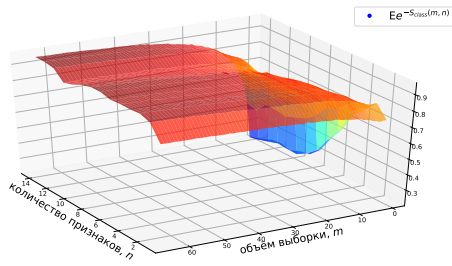
(b) Diabetes



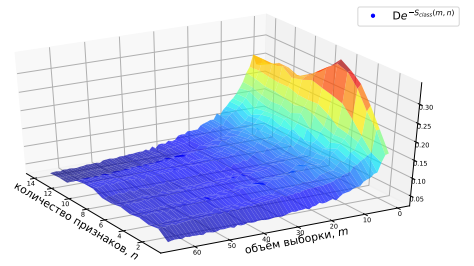
(c) Boston



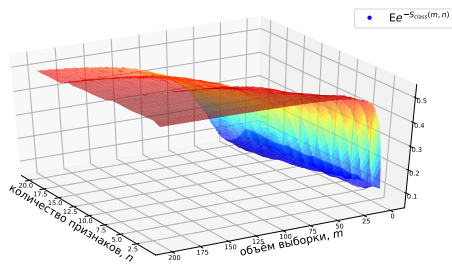
(d) Boston



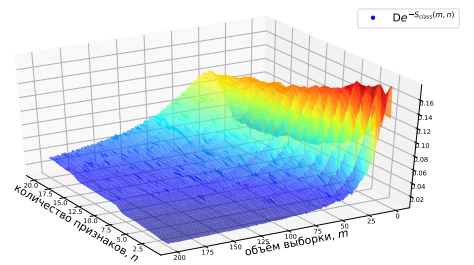
(e) Wine



(f) Wine



(g) Nba



(h) Nba

Рис. 3: Зависимость значения функции $e^{-S(m,n)}$ от объема выборки m и количества параметров n для выборок из UCI репозитория

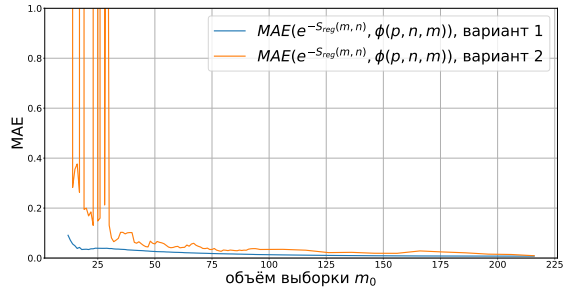
На рис. 4. представлены графики качества аппроксимации функции $\hat{l}(m)$, а также предсказание \hat{m}^* при различных m_0 для разных выборок. Так же, как и для синтетических выборок, аппроксимация функции $\hat{l}(m)$ в варианте с неполной информацией даёт качество хуже, чем в варианте с полной информацией. Также качество аппроксимации в варианте с неполной информацией сильно шумит при малых значениях m_0 .

При $m_0 \rightarrow m$ аппроксимация в варианте с неполной информацией стремится к аппроксимации в варианте с полной информацией, поэтому при больших m_0 предсказания \hat{m}^* для этих двух вариантов практически совпадают.

Для выборок задачи регрессии получилось построить адекватное предсказание достаточного объёма выборки \hat{m}^* при $m_0 < m^*$.

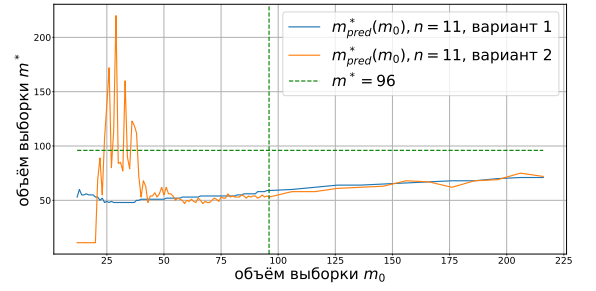
Для выборок задачи классификации при $m_0 < m^*$ построить адекватное предсказание \hat{m}^* не получилось. Это может быть связано с тем, что достаточный объём m^* для задачи классификации сильно меньше, чем для задачи регрессии, а при меньших m_0 аппроксимация функции $\hat{l}(m)$ хуже, чем при больших m_0 .

Аппроксимация $e^{-S(m,n)}$

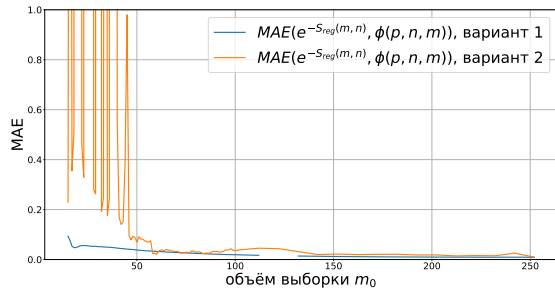


(a) Diabetees

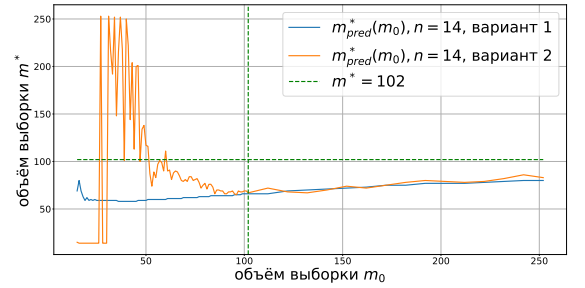
Предсказание m^*



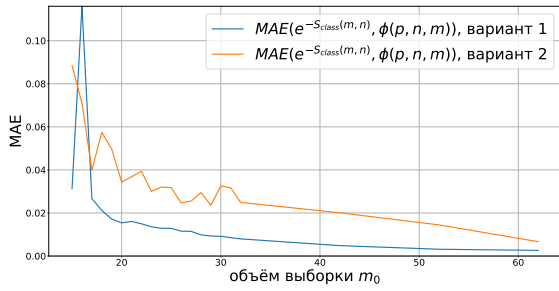
(b) Diabetees



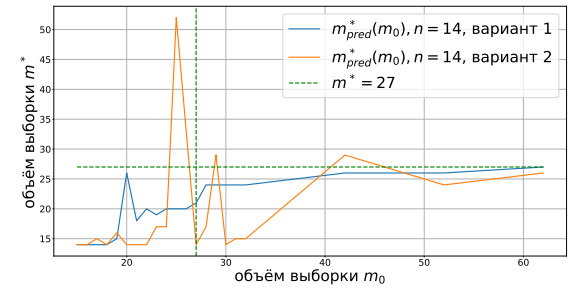
(c) Boston



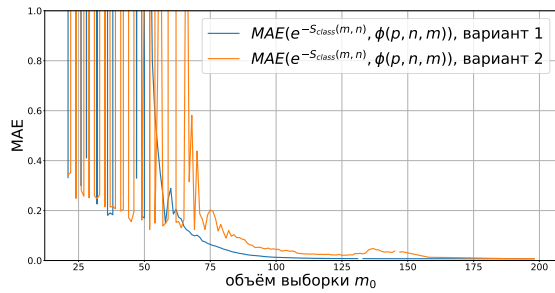
(d) Boston



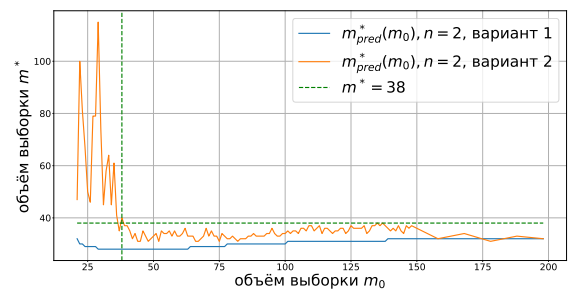
(e) Wine



(f) Wine



(g) Nba



(h) Nba

Рис. 4: Качество предсказания $e^{-S(m,n)}$ и m^* в зависимости от объема обучающей выборки m_0 для выборок из UCI репозитория

4. Заключение

Построен алгоритм раннего прогнозирования достаточного объёма выборки для обобщённой линейной модели. Для некоторых выборок задача раннего прогнозирования была успешно решена. Данный алгоритм можно улучшить путём использования в построении модели аппроксимации $\phi(m) \sim \hat{l}(m)$ свойств оценки вектора параметров $\hat{\mathbf{w}}$ обобщенной линейной модели.

Список литературы

- [1] S. G. Self and R. H. Mauritsen Power/sample size calculations for generalized linear models // Biometrics, 1988. Vol. 44. P. 79-86.
- [2] G. Shieh On power and sample size calculations for likelihood ratio tests in generalized linear models // Biometrics, 2000. Vol. 56. P. 1192-1196.
- [3] 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.
- [4] 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.
- [5] Rosa L Figuerola, Qing Zeng-Treitler, Sasikiran Kandula and Long H Ngo Predicting sample size required for classification performance // BMC Medical Informatics and Decision Making.
- [6] Kevin K. Dobbin, Yingdong Zhao, and Richard M. Simon How Large a Training Set is Needed to Develop a Classifier for Microarray Data? // American Association for Cancer Research, 2008.