

Исследование логистической регрессии

Петров Егор

МФТИ

2 ноября 2024 г.

- 1 Постановка задачи
 - LogLoss
 - Выпуклость
 - Оптимизация
- 2 Реализация и исследование
 - Распределение Бернулли
 - Сравнение методов оптимизации
 - Зависимость от гиперпараметров
- 3 Визуализация предсказаний
 - Визуализация
- 4 Итог

Постановка задачи

- Постановка задачи логистической регрессии

Итак, ставим задачу бинарной классификации $\{x_i, y_i\}_{i=1}^n$,
 $x_i \in \mathbb{R}^d, y_i \in \{0, 1\}$

Сама модель логистической регрессии имеет вид:

$$\hat{y}_i = \begin{cases} 0, \text{ если } \sigma(\theta^T x) < threshold \\ 1, \text{ иначе} \end{cases} \quad (1)$$

В таком случае, мы пользуемся предположением, о том, что сигмоиду $\sigma(\theta^T x) = \frac{1}{1 + \exp(-\theta^T x)}$ можно интерпретировать как вероятность

Теперь займемся выводом оптимизируемого функционала ошибки. Поскольку $y_i \in \{0, 1\}$, то будем интерпретировать задачу как поиск оптимального параметра для распределения $Bern(p)$, то есть $y_i \sim Bern(p)$.

Теперь перепишем условную вероятность в виде:

$$\mathbb{P}_p(y_i) = p^{y_i} \cdot (1 - p)^{1-y_i} \quad (2)$$

Тогда правдоподобие примет вид:

$$\mathbb{L}_{y_i}(p) = p^{\sum_{i=1}^n y_i} \cdot (1 - p)^{n - \sum_{i=1}^n y_i} \quad (3)$$

Рассмотрим логарифм правдоподобия (здесь стоит заметить, что логарифм монотонная возрастающая функция, а значит с точки зрения оптимизации данный переход корректен)

$$l_{y_i}(p) = \left(\sum_{i=1}^n y_i \right) \log(p) + \left(n - \sum_{i=1}^n y_i \right) \log(1 - p) \quad (4)$$

Теперь заменим p на нашу оценку истинной вероятности $\sigma(\theta^T x_i)$ и и записываем под одной суммой

$$l_{y_i}(\theta) = \sum_{i=1}^n y_i \log(\sigma(\theta^T x_i)) + (1 - y_i) \log(1 - \sigma(\theta^T x_i)) \rightarrow \max_{\theta} \quad (5)$$

Тогда сама функция $\text{LogLoss}(\theta)$ приме вид (заменяем для минимизации знак в предыдущем выражении)

$$\text{LogLoss}(\theta) = - \left(\sum_{i=1}^n y_i \log(\sigma(\theta^T x_i)) + (1 - y_i) \log(1 - \sigma(\theta^T x_i)) \right) \rightarrow \min_{\theta} \quad (6)$$

Таким образом, получили оптимизируемый функционал
Теперь убедимся, что для нахождения оптимума полученного выражения можно пользоваться итерационными методами, основанными на градиентном спуске, то есть докажем выпуклость полученного функционала

Выпуклость

Для доказательства этого факта воспользуемся критерием выпуклости второго порядка, который гласит

$$f(x) - \text{выпукла} \Leftrightarrow \text{dom}(f) - \text{выпуклое и } \nabla^2 f \succcurlyeq 0 \quad (7)$$

В данном случае, выпуклость множества $\text{dom}(f)$ очевидна, поскольку следует из вида нашей функции, так как здесь $\text{dom}(f) = \mathbb{R}^d$, что по определению является выпуклым множеством (поскольку значения LogLoss по модулю близкие к $\pm\infty$ достигаются при аналогичных значениях $\theta \rightarrow \pm\infty$, которые входят только в расширенное пространство $\overline{\mathbb{R}}$)

Теперь остановимся на положительной полуопределенности гессиана, для начала найдем его

Градиент:

$$\nabla_{\theta} \text{LogLoss}(\theta) = -\nabla_{\theta} \left(\sum_{i=1}^n y_i \log(\sigma(\theta^T x_i)) + (1 - y_i) \log(1 - \sigma(\theta^T x_i)) \right) =$$

$$= -\sum_{i=1}^n \left(y_i \nabla_{\theta} \log(\sigma(\theta^T x_i)) + (1 - y_i) \nabla_{\theta} \log(1 - \sigma(\theta^T x_i)) \right) \quad (8)$$

Далее следует заметить, что $\sigma'(x) = \sigma(x)(1 - \sigma(x))$ $\nabla_{\theta} \text{LogLoss}(\theta) =$

$$= -\sum_{i=1}^n \left(y_i \frac{\sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i)) \cdot x_i}{\sigma(\theta^T x_i)} - (1 - y_i) \frac{\sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i)) \cdot x_i}{1 - \sigma(\theta^T x_i)} \right) =$$

$$= -\sum_{i=1}^n \left(y_i(1 - \sigma(\theta^T x_i))x_i - (1 - y_i)\sigma(\theta^T x_i) \cdot x_i \right) =$$

$$= -\sum_{i=1}^n \left(y_i - \sigma(\theta^T x_i) \right) x_i = \sum_{i=1}^n \left(\sigma(\theta^T x_i) - y_i \right) x_i$$

Что можно несложно переписать в матричном виде:

$$\nabla_{\theta} \text{LogLoss}(\theta) = X^T (S(\theta) - Y) \quad (9)$$

где $X_i = x_i$, $S(\theta)_i = \sigma(\theta^T x_i)$, $Y_i = y_i$

Теперь вычислим непосредственно гессиан:

$$\nabla_{\theta}^2 \text{LogLoss}(\theta) = \sum_{i=1}^n (\sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i))) x_i x_i^T \quad (10)$$

Что также несложно переписывается в матричном виде:

$$\nabla_{\theta}^2 \text{LogLoss}(\theta) = X^T V(\theta) X \quad (11)$$

где $V(\theta) = \text{diag}(\sigma(\theta^T x_i)(1 - \sigma(\theta^T x_i)))$

Заметим, что $V(\theta) \succcurlyeq 0$, поскольку матрица диагональна и при этом на самой диагонали стоят положительные элементы, а поскольку эта матрица умножается на X^T и X , то и в произведении получим положительно полуопределенную матрицу

Таким образом, доказали, что

$$\nabla_{\theta}^2 \text{LogLoss}(\theta) \succcurlyeq 0 \quad (12)$$

А значит, сама функция $\text{LogLoss}(\theta)$ - выпуклая

Теперь выпишем итеративные способы нахождения оптимального параметра θ (нетрудно показать, что при добавлении регуляризации формулы примут следующий вид):

Итеративная формула для GD:

$$\theta_{k+1} = \theta_k - \eta \cdot (-X^T(Y - S(\theta_k))) + 2\eta\lambda\theta_k \quad (13)$$

Итеративная формула для SGD:

$$\theta_{k+1} = \theta_k - \eta \frac{n}{|I|} \sum_{i \in I} (-X_i(Y_i - S_i(\theta_k))) + 2\eta\lambda\theta_k \quad (14)$$

Согласно формулам полученным ранее матрица Гессе имеет вид:

$$\nabla^2 F(\theta) = X^T V(\theta) X + 2\lambda E, \text{ где} \quad (15)$$

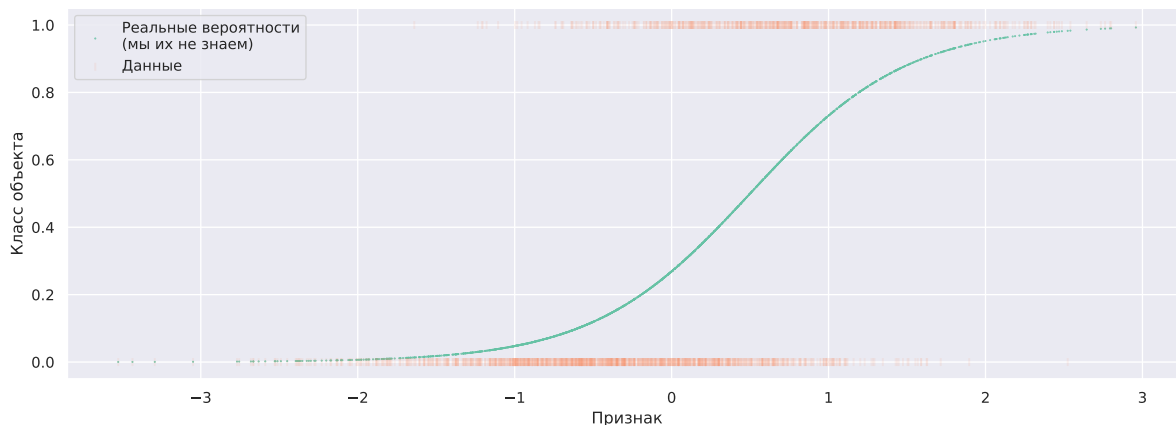
$V(\theta) = \text{diag}(\sigma(X_i^T \theta)(1 - \sigma(X_i^T \theta)))$, E —единичная матрица Тогда итеративная формула для IRLS:

$$\theta_{k+1} = \theta_k - (X^T V(\theta_k) X + 2\lambda E)^{-1} \cdot (-X^T (Y - S(\theta_k)) + 2\lambda \theta_k), \text{ где} \quad (16)$$

$V(\theta) = \text{diag}(\sigma(X_i^T \theta)(1 - \sigma(X_i^T \theta)))$, E —единичная матрица

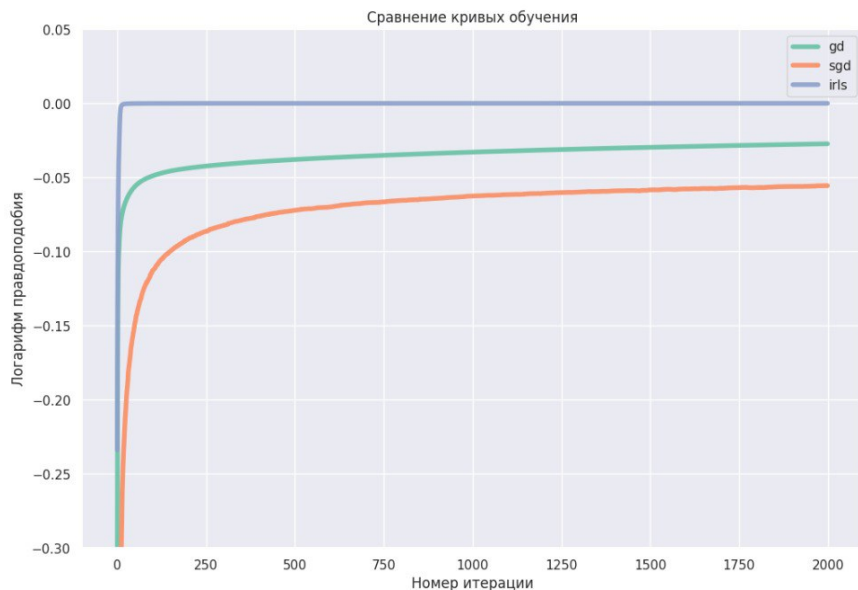
Распределение Бернулли

Для начала изобразим выборку из распределения Бернулли (Данные), которые по оси X распределены нормально, а также отрисуем вероятности принадлежности классам 0 и 1, которые мы получим при аппроксимации вероятности сигмоидой



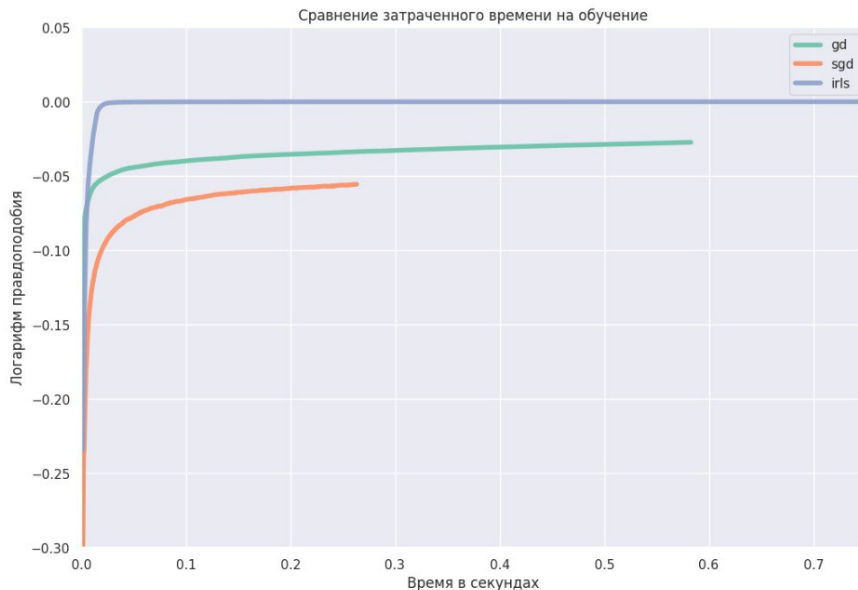
Сравнение методов оптимизации

Реализацию логистической регрессии можно посмотреть в прикрепленном файле. В ней для оптимизации функционала ошибки дано на выбор три функции, описанные выше. Сравним их



Сравнение методов оптимизации

Теперь посмотрим на затраченное время для каждого из методов



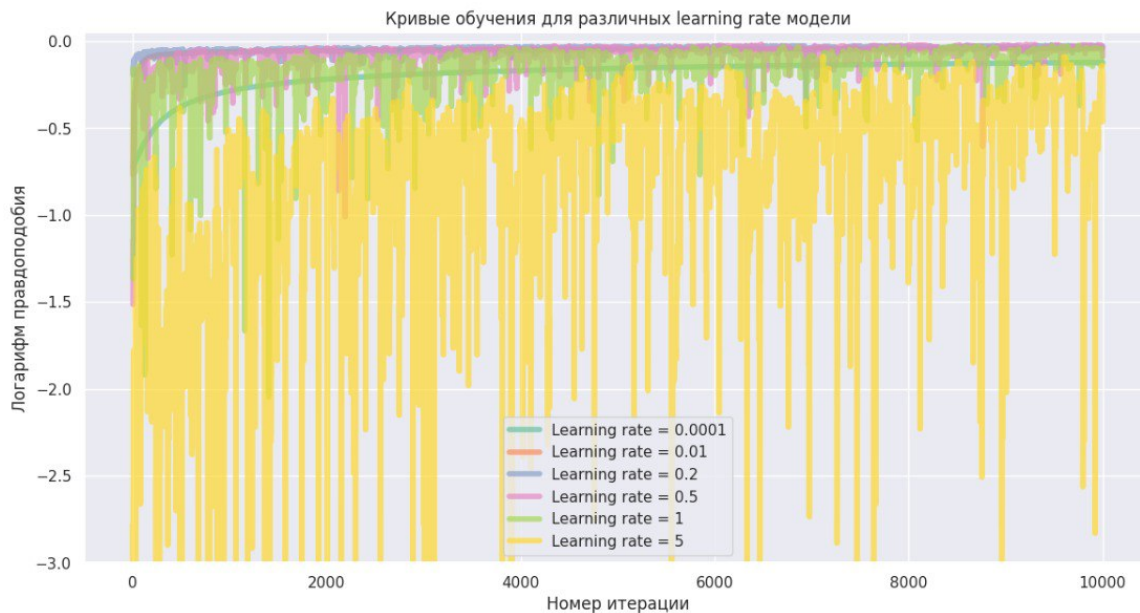
Зависимость от гиперпараметров

В данной секции исследуем зависимость качества логистической регрессии (по метрике Accuracy и методу оптимизации gd) от learning rate и коэффициента регуляризации
Начнем с learning rate



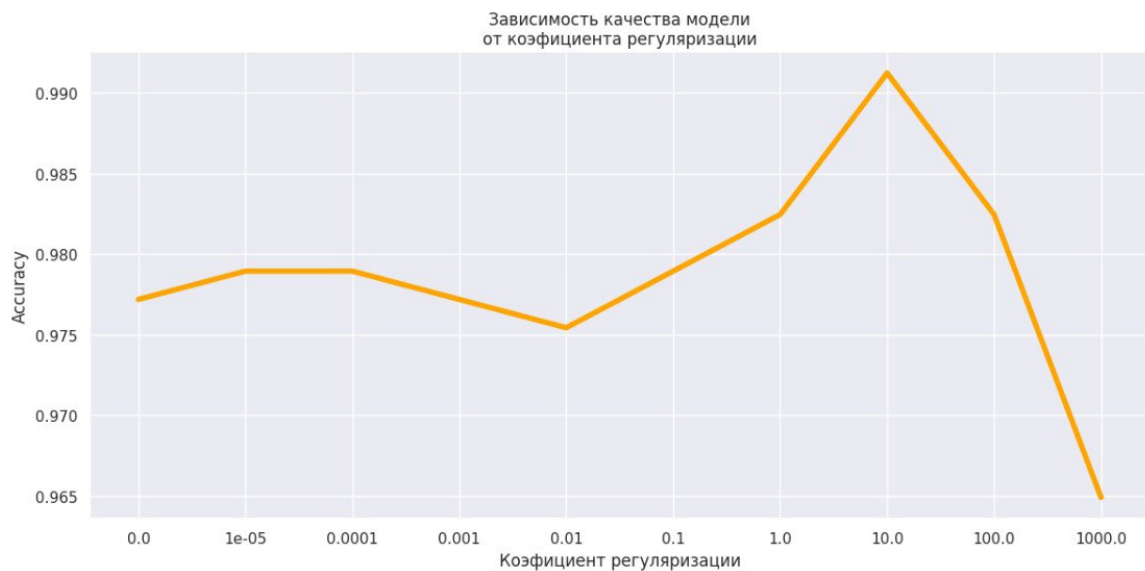
Зависимость от гиперпараметров

В данном случае важно отразить и различия в обучении для перебираемых значений lr



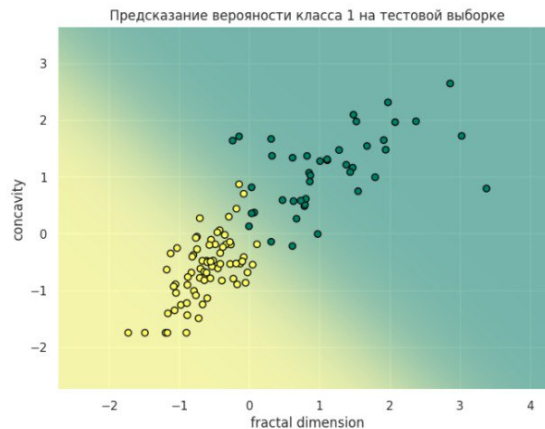
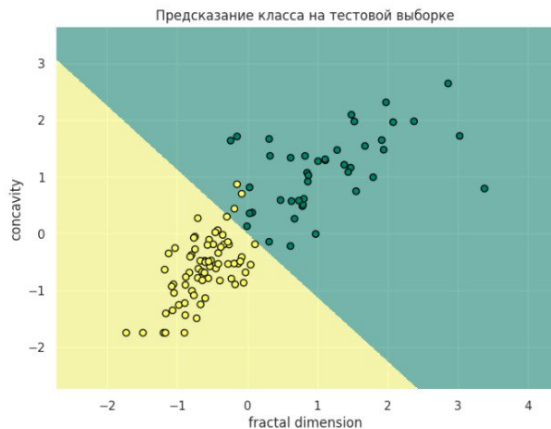
Зависимость от гиперпараметров

Теперь изучим модель при изменении параметра регуляризации



Визуализация предсказаний

Теперь непосредственно изобразим предсказания нашей модели, визуализировав как непосредственно предсказания, так и распределения вероятностей



Особенности логистической регрессии

- Логистическая регрессия может отлично справляться с задачей бинарной классификации, когда классы хорошо разделимы при помощи гиперплоскости, поскольку в её основе лежит линейная модель
- Важное предположение, которые мы делаем - вероятность хорошо аппроксимируется сигмой (само это утверждение конечно требует доказательство, которое строится, основываясь на Байесовском подходе)
- Нами были выведен функционал ошибки для решения задачи бинарной классификации, а также доказана его выпуклость, позволяющая решать задачу градиентными методами