

Машинное обучение, DS-поток

Домашнее задание 2

Правила:

- Дедлайн **28 февраля 10:00**. После дедлайна работы не принимаются кроме случаев наличия уважительной причины.
- Выполненную работу нужно отправить на почту `mipt.stats@yandex.ru`, указав тему письма "[ml] Фамилия Имя - задание 2". Квадратные скобки обязательны. Если письмо дошло, придет ответ от автоответчика.
- Прислать нужно ноутбук и его pdf-версию (без архивов). Названия файлов должны быть такими: `2.N.ipynb` и `2.N.pdf`, где `N` - ваш номер из таблицы с оценками.
- Решения, размещенные на каких-либо интернет-ресурсах не принимаются. Кроме того, публикация решения в открытом доступе может быть приравнена к предоставлению возможности списать.
- Для выполнения задания используйте этот ноутбук в качестве основы, ничего не удаляя из него.
- Никакой код из данного задания при проверке запускаться не будет.

Баллы за задание:

- Задача 1 - 5 баллов
- Задача 2 - 15 баллов

In []:

```
1 %matplotlib inline
2
3 import numpy as np
4 import pandas as pd
5 import seaborn as sns
6 import scipy.stats as sps
7 import matplotlib.pyplot as plt
8
9 from sklearn.metrics import accuracy_score
10 from sklearn.datasets import load_breast_cancer
11 from sklearn.preprocessing import StandardScaler
12 from sklearn.model_selection import train_test_split
13
14 import warnings
15
16 sns.set_style("dark")
17 sns.set(font_scale=1.4)
18 warnings.filterwarnings('ignore')
```

Задача 1

Рассмотрим метод логистической регрессии. Пусть $x_i \in \mathbb{R}^d$, $Y_i \sim \text{Bern}(\mu_\theta(x_i))$.

Мы предполагаем, что $\mu_\theta(x_i) = P_\theta(Y_i = 1) = \sigma(x_i^T \theta) = \frac{e^{x_i^T \theta}}{1 + e^{x_i^T \theta}}$.

Регуляризацию в методе логистической регрессии можно задать с помощью введения априорного распределения на θ . Будем считать, что априорное распределение $\mathcal{N}(0, \alpha^{-1} I_d)$. В данном случае распределения не являются сопряженными, поэтому простым путем найти апостериорное распределение не получится. Однако, можно найти моду этого распределения. Выпишите соответствующую задачу оптимизации.

Для данной задачи:

1. Получите формулу градиентного спуска.
2. Получите формулу метода IRLS.

Задача 2

1.

Реализуйте логистическую регрессию с регуляризацией для трех вариантов поиска оценки параметров:

- обычный градиентный спуск;
- стохастический mini-batch градиентный спуск, размер батча 5-10;
- IRLS.

Для измерения времени работы **каждого** шага используйте

```
from time import time
```

Замечание. Для чистоты эксперимента время шага внутри цикла `for` нужно замерять от конца предыдущего шага до конца текущего, а не от начала текущего шага.

In []:

```
1 class LogisticRegression():
2     '''
3     Модель логистической регрессии. Имеет следующие гиперпараметры:
4
5     * alpha: параметр регуляризации.
6         Если равно 0, то регуляризация не происходит.
7     * lr: константа, на которую домножаем градиент при обучении
8     * eps: ограничение на норму невязки в случае
9         если используется критерий criterion='eps'
10    * max_iter: ограничение на кол-во итераций в случае
11        если используется критерий criterion='max_iter'
12    * method: если равно 'gd', то используется обычный градиентный спуск,
13        если равно 'sgd', то используется стохастический
14        градиентный спуск,
15        если равно 'irls', то используется метод IRLS.
16    * criterion: если равно 'eps', то используем ограничение
17        на норму невязки,
18        если равно 'max_iter', то используем ограничение
19        на количество итераций
20    * fit_intercept: указывает, следует ли добавить константу в признаки
21    * save_history: указывает, следует ли сохранять историю обучения
22    '''
23
24
25    def __init__(self, alpha=0, lr=0.5, eps=1e-3, max_iter=1e5,
26                  method='gd', criterion='max_iter',
27                  fit_intercept=True, save_history=True):
28        ''' Создает модель и инициализирует параметры '''
29
30        assert criterion in ['max_iter', 'eps'], 'выбран неправильный критерий'
31        assert method in ['gd', 'sgd', 'irls'], 'выбран неправильный метод'
32
33        self.alpha = alpha
34        self.lr = lr
35        self.eps = eps
36        self.max_iter = max_iter
37        self.criterion = criterion
38        self.method = method
39        self.fit_intercept = fit_intercept
40        self.save_history = save_history
41        self.history = [] # для хранения истории обучения
42
43
44    @staticmethod
45    def _sigmoid(x):
46        return 1 / (1 + np.exp(-x))
47
48
49    def _compute_loss(self, X, y):
50        return <...>
51
52
53    def _add_intercept(self, X):
54        '''
55        Добавляем свободный коэффициент к нашей модели.
56        Это происходит путем добавления вектора из 1 к исходной матрице.
57        '''
58
59        X_copy = np.full((X.shape[0], X.shape[1] + 1), fill_value=1)
```

```

60     X_copy[:, :-1] = X
61
62     return X_copy
63
64
65 def fit(self, X, Y):
66     '''
67     Обучает модель логистической регрессии с помощью выбранного метода,
68     пока не выполнится критерий остановки self.criterion.
69     Также, в случае self.save_history=True, добавляет в self.history
70     текущее значение оптимизируемого функционала
71     и время обновления коэффициентов.
72     '''
73
74     assert X.shape[0] == Y.shape[0]
75
76     if self.fit_intercept: # добавляем свободный коэффициент
77         X_copy = self._add_intercept(X)
78     else:
79         X_copy = X.copy()
80
81     <...>
82
83     self.coef_ = <...> # коэффициенты модели
84     self.intercept_ = <...> # свободный коэффициент
85     self.n_iter_ = <...> # произведенное число итераций
86
87     return self
88
89
90 def predict(self, X):
91     '''
92     Применяет обученную модель к данным
93     и возвращает точечное предсказание (оценку класса).
94     '''
95
96     if self.fit_intercept:
97         X_copy = self._add_intercept(X)
98     else:
99         X_copy = X.copy()
100
101     assert X_copy.shape[1] == self.weights.shape[0]
102
103     <...>
104
105     return predictions # shape = (n_test,)
106
107
108 def predict_proba(self, X):
109     ''' Применяет обученную модель к данным
110     и возвращает предсказание вероятности классов 0 и 1. '''
111
112     if self.fit_intercept:
113         X_copy = self._add_intercept(X)
114     else:
115         X_copy = X.copy()
116
117     assert X_copy.shape[1] == self.weights.shape[0]
118
119     <...>
120

```

```
121 |         return prob_predictions # shape = (n_test, 2)
```

Рассмотрим игрушечный датасет на 30 признаков `load_breast_cancer` из библиотеки `sklearn`. Это относительно простой для двуклассовой классификации датасет по диагностике рака молочной железы.

Ради интереса можно прочитать описание признаков.

In []:

```
1 dataset = load_breast_cancer()
2 dataset['DESCR'].split('\n')[11:31]
```

Разделим нашу выборку на обучающую и тестовую:

In []:

```
1 X, Y = dataset['data'], dataset['target']
2
3 X_train, X_test, Y_train, Y_test \
4     = train_test_split(X, Y, test_size=0.2, random_state=42)
5 X_train.shape, X_test.shape, Y_train.shape, Y_test.shape
```

При использовании регуляризации данные необходимо нормализовать. Воспользуемся для этого классом `StandardScaler` из библиотеки `sklearn`.

In []:

```
1 scaler = StandardScaler()
```

2. Теперь обучите три модели логистические регрессии без регуляризации с помощью методов

- обычный градиентный спуск;
- стохастический mini-batch градиентный спуск;
- IRLS

Постройте график, на котором нанесите три кривые обучения, каждая из которых отображает зависимость оптимизируемого функционала от номера итерации метода. Функционал должен быть одинаковый для всех моделей, то есть без минусов. Нарисуйте также график зависимости этого функционала от времени работы метода.

Для чистоты эксперимента желательно не запускать в момент обучения другие задачи и провести обучение несколько раз, усреднив результаты.

Напоминание: все графики должны быть информативны, с подписанными осями и т.д.

Сделайте выводы. Что будет, если при обучении на очень большой по количеству элементов датасете?

In []:

```
1
```

3. Сравните два реализованных критерия остановки по количеству проведенных итераций: евклидова норма разности текущего и нового векторов весов стала меньше, чем $1e-4$ и ограничение на число

итераций (например, 10000). Используйте градиентный спуск.

In []:

1

4. Рассмотрите как влияет размер шага (`learning rate`) на качество модели. Обучите каждую модель одинаковое число итераций (например, 10000), а затем посчитайте качество. Воспользуйтесь ограничением на число итераций в качестве критерия остановки, так как для больших `learning rate` у вас может не сойтись модель. Используйте стохастический градиентный спуск. Сделайте выводы.

In []:

```
1 lrs = [0.01, 0.1, 0.2, 0.3, 0.5, 0.7, 1, 2, 5, 10]
2
```

Постройте кривые обучения для различных `learning rate` . Не обязательно рассматривать все `learning rate` из предыдущего задания, так как их слишком много, и график будет нагроможден. Возьмите около половины из них. Какой `learning rate` лучше выбрать? Чем плохи маленькие и большие `learning rate` ?

In []:

1

5. Рассмотрите несколько моделей, в которых установите не менее 5-ти различных коэффициентов регуляризации, а также модель без регуляризатора. Сравните, влияет ли наличие регуляризации на скорость сходимости и качество, сделайте выводы. Под качеством подразумевается значение метрики, рассмотренных на семинаре.

In []:

1

6. Выберите произвольные два признака, в пространстве которых визуализируйте предсказания вероятностей класса 1 для модели, которая показала наилучшее качество на предыдущем шаге.

In []:

1