

## Задание 2. Метрические алгоритмы классификации

Практикум 317 группы, 2014

Начало выполнения задания: 14 октября 2014 года.

Срок сдачи: **4 ноября 2014 года, 23:59.**

Максимальный балл: 5.0 (плюс бонусные баллы).

Текст задания обновлён 27 октября 2014 года.

## Содержание

|          |                       |          |
|----------|-----------------------|----------|
| <b>1</b> | <b>Задание</b>        | <b>1</b> |
| <b>2</b> | <b>Бонусные баллы</b> | <b>3</b> |

## 1 Задание

Данное задание направлено на ознакомление с метрическими алгоритмами классификации, а также методами работы с текстовой информацией.

Сдача задания осуществляется через репозиторий SVN. Задание должно быть выполнено в IPython Notebook, формат сдачи — .ipynb и .ру файлы, а также опционально отчёт в формате pdf. За отсутствие анализа экспериментов и выводов балл будет снижаться!

1. Загрузить датасет MNIST при помощи функции `sklearn.datasets.fetch_mldata("MNIST original")`.
2. Разбить датасет на обучающую выборку (первые 60 тыс. объектов) и тестовую выборку (10 тыс. последних объектов). Ответы на тестовой выборке не следует использовать ни в каких экспериментах, кроме финального.
3. Визуализировать по 5 случайных объектов из каждого из 10 классов. Воспользуйтесь методами `np.reshape`, `pyplot.subplot`, и `pyplot.imshow` с параметром `cmap="Greys"`. Также можно убрать оси координат при помощи команды `pyplot.axis("off")`.
4. Исследовать, каким точным алгоритмом поиска ближайших соседей следует пользоваться в различных ситуациях. Будем искать 5 ближайших соседей в обучающей выборке для тестовой выборки (при этом ответы на тестовой выборке не используются!). Метрика евклидова. Число признаков: 10, 20, 100. Подмножество признаков выбирается один раз, случайно. Алгоритмы поиска ближайших соседей:
  - (a) Собственная реализация на основе кода подсчёта евклидова расстояния между двумя множествами точек из задания №1
  - (b) `sklearn.neighbors.NearestNeighbors(algorithm='brute')`
  - (c) `sklearn.neighbors.NearestNeighbors(algorithm='kd_tree')`

**Замечание 1.** При поиске  $k$  ближайших соседей некоторые методы строят в памяти матрицу попарных расстояний обучающей выборки и тестовой выборки. Рекомендуем написать функцию, которая ищет ближайших соседей блоками, то есть делает запросы ближайших соседей для первых  $N$  тестовых объектов, затем для следующих  $N$ , и так далее, и в конце объединяет полученные результаты.

**Замечание 2.** Для оценки времени долго работающих функций можно пользоваться либо командой `time.clock()`, либо magic-командой `%time`, которая запускает код лишь один раз.

5. Реализовать генерацию индексов обучающей и валидационной выборки для кросс-валидации с  $p$  фолдами. Не разрешается использовать модуль `sklearn.cross_validation`.
6. Пусть дана обучающая и валидационная выборка. Реализовать оценку точности метода  $k$  ближайших соседей по валидационной выборке для нескольких параметрах  $k$ :  $[k_1, \dots, k_n]$ ,  $k_1 < k_2 < \dots < k_n$ . Сложность алгоритма для одного объекта из валидационной выборки должна иметь порядок  $O(k_n)$ . Нельзя использовать класс `sklearn.neighbors.KNeighborsClassifier`, можно пользоваться готовыми реализациями поиска ближайших соседей.

7. Оценить по кросс-валидации с 3 фолдами точность (долю правильно предсказанных ответов) и скорость метода k ближайших соседей в зависимости от следующих факторов:
  - (a) k от 1 до 10.
  - (b) Используется евклидова или косинусная метрика.
8. Реализовать взвешенный метод k ближайших соседей, где голос объекта равен  $1/(distance + \epsilon)$ , где  $\epsilon$  — малое число. Сравнить с методом без весов при тех же фолдах и тех же параметрах.
9. Написать модуль `contest.py`, содержащий функцию `get_mnist_predictions(train, train_answer, test)`. Функция принимает обучающую выборку `train`, ответы на обучающей выборке `train_answer` и тестовую выборку `test`, и должна вернуть предсказания на тестовой выборке. Разрешается модифицировать метод k ближайших соседей, использовать другие классификаторы, модифицировать признаки. Можно пользоваться готовыми реализациями алгоритмов, но, в случае использования необычных пакетов, просьба приложить инструкцию по установке. Функция должна работать не дольше 5 минут на компьютере преподавателя, не должна производить перебор параметров и не может использовать вспомогательные данные (в том числе нельзя загружать уже обученный классификатор; обучение должно производиться внутри функции). После сдачи задания будет проведено тестирование методов на скрытом разбиении на обучение и тест. Авторы пяти лучших методов получают +0.5 балла.
10. Применить функцию `get_mnist_predictions` к исходной обучающей и тестовой выборке. Подсчитать точность. Сравнить с точностью по кросс-валидации. Сравнить с указанной в интернете точностью лучших алгоритмов на данной выборке.
11. Визуализировать несколько объектов из тестовой выборки, на которых были допущены ошибки. Проанализировать и указать их общие черты. Построить и проанализировать матрицу ошибок (`confusion matrix`). Можно воспользоваться функцией `sklearn.metrics.confusion_matrix`.
12. Загрузить обучающую выборку датасета `20 newsgroups` при помощи метода `sklearn.datasets.fetch_20newsgroups`. Убрать все заголовки, подписи и цитаты, используя аргумент `remove`.
13. Перевести во всех документах все буквы в нижний регистр. Заменить во всех документах символы, не являющиеся буквами и цифрами, на пробелы. Полезные функции: `str.lower`, `str.isalnum`.
14. Разбить каждый документ на термы по пробельным символам (пробелы, переносы строки). Полезная функция: `str.split`.
15. Преобразовать датасет в разреженную матрицу `scipy.sparse.csr_matrix`, где значение `x` в позиции `(i, j)` означает, что в документе `i` слово `j` встретилось `x` раз. Необходимо воспользоваться наиболее эффективным конструктором `csr_matrix((data, indices, indptr), shape=(M, N))`.

**Замечание.** Не забудьте указать параметр `shape`, так как число используемых термов в тестовой выборке может отличаться от числа термов в обучении.
16. Произвести `tf-idf` преобразование датасета при помощи `sklearn.feature_extraction.text.TfidfTransformer`. Используйте параметры по умолчанию.
17. Оценить точность (долю правильно предсказанных ответов) и скорость метода k ближайших соседей при помощи кросс-валидации с 3 фолдами. Точность для всех значений k должна оцениваться в один проход по валидационной выборке. Исследуйте на одних и тех же фолдах влияние следующих факторов:
  - (a) k от 1 до 10.
  - (b) Используется евклидова или косинусная метрика.
  - (c) Используется ли преобразование `tf-idf`.
  - (d) Используются взвешенный или невзвешенный метод k ближайших соседей.
18. Загрузите тестовую выборку (параметр `subset` метода `fetch_20newsgroups`). Примените лучший алгоритм к тестовой выборке. Сравните точность с полученной по кросс-валидации.
19. Вывести несколько документов из тестовой выборки, на которых были допущены ошибки. Проанализировать их. Построить и проанализировать матрицу ошибок (`confusion matrix`). Можно воспользоваться функцией `sklearn.metrics.confusion_matrix`.
20. Сделать выводы: в каких случаях следует пользоваться каким алгоритмом, какие параметры важно оптимизировать, какие ошибки допускают алгоритмы.

## 2 Бонусные баллы

- +0.25 балла. Качественное проведение дополнительного (не пересекающегося с основным заданием) мини-исследования по теме метрических алгоритмов: формулируется изучаемый вопрос, ставятся эксперименты, позволяющие на него ответить, делаются выводы.
- +0.5 балла. После сдачи задания будет составлен рейтинг решений задачи MNIST по точности на скрытом разбиении на обучение и контроль. Авторы пяти лучших решений получат +0.5 балла.