



СОФИЙСКИ УНИВЕРСИТЕТ „СВ.
КЛИМЕНТ ОХРИДСКИ“
ФАКУЛТЕТ ПО МАТЕМАТИКА И
ИНФОРМАТИКА

КУРСОВ ПРОЕКТ
ПО СИСТЕМИ, ОСНОВАНИ НА
ЗНАНИЯ

Тема:

Разпознаване на авторство

Студент:

Симеон Емилов Христов, Група 3, ФН: 71845

София, януари 2021 г.

1. Формулировка на задачата

Да се реализира класификатор, който разпознава текстове на Иван Вазов от текстове на Йордан Йовков.

Класификаторът е написан на Python. Съставен е от 4 различни файла, разработени в Google Colaboratory и съответно с разширение “.ipynb”. За правилно тестване е необходимо предварително качване на текстовите материали, изброени в т. 3.

2. Използвани алгоритми

Използва се вид машинно самообучение с учител - самообучение чрез запомняне. Реализиран е методът на логистичната регресия. За характеристики на входните данни се създава таблица от вида “честота на думата - обратна честота на документа”. За оценяване работата на логистичната регресия се използват разнообразни методи, включ. матрица за разбиране на объркването.

Най-често посещаваните източници са документациите на Python, и библиотеките Pandas, Sklearn, Numpy, Matplotlib и Pickle.

3. Описание на програмната реализация

Проектът се състои от 4 файла, всеки представляващ отделна стъпка в процеса на машинното самообучение чрез запомняне:

1. 1_data_cleaning.ipynb (цел: създаване на corpus и tfidf vectorizer):
 - a. извличане на текстовете, използвани за обучаващи примери.
 - i. текстовете са: “Под игото” (Иван Вазов), “Епопея на забравените” (Иван Вазов), “Чифликът край границата” (Йордан Йовков), “Приключенията на Гороломов” (Йордан Йовков), “Старопланински легенди” (Йордан Йовков), “Последна радост” (Йордан Йовков), “Вечери в Антимовския храм” (Йордан Йовков).
 - b. формиране на “тяло” - таблица с две колони, представляващи съответно текст и авторът му.
 - c. “изчистване” на отделните текстове:
 - i. заместване на главните букви със съответните им малки.
 - ii. премахване на пунктуация.
 - iii. премахване на думи, които не се състоят от български букви и последователности от символи, в които има такива, които не са букви.

- iv. премахване на междуметия и често използвани словосъчетания.
 - v. заместване на думи с коренната им такава.
 - d. създаване на биграми.
 - e. създаване на таблица "честота на думата - обратна честота на документа".
- 2. 2_exploratory_data_analysis.ipynb (цел: извличане на закономерности от данните):
 - a. генериране на таблица с обща информация за всички колони, които съдържат числа.
 - b. намиране отговор на въпросите: "Коя е "най-тежката" дума? Какво е нейното тегло? Кой автор я е използвал? Коя дума има най-висока средна тежест? Колко е тя? Кой автор я е използвал?"
 - c. създаване на хистограми, показващи взаимоотношенията между различните думи.
- 3. 3_apply_techniques.ipynb (цел: чрез използване на логистична регресия да се създаде модел):
 - a. проверка за най-подходящ модел измежду LogisticRegression, KNeighborsClassifier, GaussianNB, MultinomialNB чрез използване на k итеративни разделяния.
 - b. чрез използване на модел, реализиращ логистична регресия да се намерят думите, отличаващи най-много двамата автори.
- 4. 4_testing_ground.ipynb (цел: да се провери и оцени работата на модела чрез използване на 150 случайно подбрани части от текстове):
 - a. формиране на случайна извадка от всички текстове.
 - b. за всеки текст от извадката да се вземе случаен интервал от символи.
 - c. прилагане на модела върху всеки откъс.
 - d. оценка на метода:
 - i. намиране на процент коректни класификации;
 - ii. генериране на матрица за разбиране на объркването;
 - iii. генериране на метрики, произхождащи от матрицата за разбиране на объркването;
 - iv. построяване на хистограма, показваща честотите на вероятностите за класификации на Иван Вазов;
 - v. построяване на графика, показваща зависимостта между реалните правилни класификации и нереалните правилни класификации.

4. Примери, илюстриращи работата на програмната система

Началото на файл със 150 случайно генерирани входи (параграфите са разделени със символите “@@@”):

[illegible]

Изходи (True - коректно класифициран, 0 - некоректно класифициран):

[illegible]

5. Литература

- <https://www.analyticsvidhya.com/blog/2018/02/the-different-methods-deal-text-data-predictive-python/> : дата на първо посещение (26.12.2020г.)

- Python 3 documentation: <https://docs.python.org/3/tutorial/index.html> дата на първо посещение (26.12.2020г.)
- Sklearn documentation: <https://sklearn.org/documentation.html> дата на първо посещение (26.12.2020г.)
- Pandas documentation: <https://pandas.pydata.org/pandas-docs/stable/reference/index.html> дата на първо посещение (26.12.2020г.)
- Numpy documentation: <https://numpy.org/doc/1.20/reference/index.html> дата на първо посещение (26.12.2020г.)
- Matplotlib documentation: <https://matplotlib.org/contents.html> дата на първо посещение (26.12.2020г.)
- Pickle documentation: <http://docs.picklesdoc.com/en/latest/> дата на първо посещение (26.12.2020г.)
- Много полезни Youtube канали:
 - [Ken Jee](#) дата на първо посещение (26.12.2020г.);
 - [Data Professor](#) дата на първо посещение (26.12.2020г.).