

# HW #06: Spark RDD

**Deadline: 29.07.2019, 08:00**

---

1. Описание задания.	1
2. Критерии оценивания.	1
3. Описание данных	2
4.1 Задача #1: народные биграммы.	2
4.2 Задача #2: коллокации.	3
5. Сроки сдачи и правила оформления задания.	4
6. Дорешка.	5

---

## 1. Описание задания.

В данном ДЗ нужно решить **2 задачи**. Задачи общие для всех. Решение надо выполнить с помощью Apache Spark, можно использовать только RDD API.

## 2. Критерии оценивания.

Балл за задачу складывается из:

- 60% - правильное решение задачи
- 20% - поддерживаемость и читаемость кода (Clean Code, см. например [Google Python Style Guide](#))
- 20% - эффективность решения

Штрафы:

- 10% за несоответствие правилам оформления задания;
- 30% за просрочку дедлайн;



Веса задач:

1. 40%
2. 60%

## 3. Описание данных

### 3.1 Дамп Википедии

en\_articles\_part:

- Путь на кластере: полный<sup>1</sup> датасет - /data/wiki/en\_articles\_part
- Формат: текст
- В каждой строке находятся следующие поля, разделенные знаком табуляции:
  1. INT - id статьи,
  2. STRING - текст статьи,

*Пример:*

```
12      Anarchism      Anarchism is often defined as a political
philosophy which holds the state to be undesirable, unnecessary, or
harmful.
```

### 3.2 Стоп-слова

stop\_words\_en:

- Путь на кластере: /data/wiki/stop\_words\_en-хро6.txt
- Формат: одно стоп-слово на строку

*Пример:*

```
...
wherein
whereupon
wherever
...
```

## 4.1 Задача #1<sup>2</sup>: народные биграммы.

Найдите все пары двух последовательных слов (биграмм), где первое слово:

```
narodnaya
```

---

<sup>1</sup> Да, здесь нет ошибки, работаем на части данных, чтобы побыстрее познакомиться со Spark RDD

<sup>2</sup> Внутренний ID задачи (для проверяющих) - 517



Для каждой пары подсчитайте количество вхождений в тексте статей Википедии. Выведите все пары с их частотой вхождений в лексикографическом порядке. Формат вывода:

```
word_pair <tab> count
```

Условия:

- очистить тексты от знаков пунктуации (см. `re.sub`)
- привести все слова к нижнему регистру;
- слова в паре объединить символом нижнего подчеркивания “\_”;
- отсортировать слова в выводе по алфавиту;

Пример вывода:

```
...  
crazy_zoo 42  
red_apple 100500  
...
```

## 4.2 Задача #2<sup>3</sup>: коллокации.

Коллокация - это комбинации слов, которые часто встречаются вместе. Например, «high school» или «Roman Empire». Чтобы определить является ли пара слов коллокацией, можно воспользоваться метрикой NPMI - нормализованная точечная взаимная информация.

Чтобы рассчитать NPMI введем несколько определений:

1.  $P(a)$  - вероятность увидеть слово “a” в датасете.  
$$P(a) = \text{num\_of\_occurrences\_of\_word\_} "a" / \text{total\_number\_of\_words}$$
  
`total_number_of_words` - общее количество слов в тексте
2.  $P(ab)$  - вероятность увидеть пару слов “a” и “b”, идущих подряд.  
$$P(ab) = \text{num\_of\_occurrences\_of\_pair\_} "ab" / \text{total\_number\_of\_word\_pairs}$$
  
`total_number_of_word_pairs` - общее количество пар
3.  $PMI(a,b) = \ln( P(ab) / [P(a) \times P(b)] )$
4.  $NPMI(a,b) = PMI(a,b) / -\ln(P(ab))$  - величина PMI нормализованная в диапазон  $[-1, 1]$ ;

Примеры и комментарии:

---

<sup>3</sup> Внутренний ID задачи (для проверяющих) - 518



- значение NPMI равное “-1” будет означать, что пара слов никогда не встречается в датасете. Например такие пары как “**green idea**” или “**sleeps furiously**” никогда не встречаются вместе, поэтому  $P(ab) = 0$ , следовательно  $PMI(a,b) = -\inf$ ,  $NPMI = -1$ ;
- значение NPMI равное “0” будет означать, что слова в паре встречается абсолютно независимо друг от друга. Рассмотрим пример “**the doors**”: “the” может встретиться рядом с любым словом. Таким образом,  $P(ab) = P(a) \times P(b)$  и  $PMI(a,b) = \ln(1) = 0$ ,  $NPMI = 0$ .
- значение NPMI равное “1” будет означать, что это идеальная коллокация. Рассмотрим пример “**the doors**”: “the” может встретиться рядом с любым словом. Таким образом,  $P(ab) = P(a) \times P(b)$  и  $PMI(a,b) = \ln 1 = 0$ ,  $NPMI = 0$ . Предположим, что “**Roman Empire**” - это уникальная комбинация, и за каждым появлением “Roman” следует “Empire”, и, наоборот, каждому появлению “Empire” предшествует “Roman”. В этом случае  $P(ab) = P(a) = P(b)$ , поэтому  $PMI(a,b) = -\ln(P(a)) = -\ln(P(b))$ , следовательно  $NPMI = 1$ .

Условия:

- найти самые популярные коллокации в Википедии;
- очистить тексты от знаков пунктуации (см. `re.sub`)
- привести все слова к нижнему регистру;
- удалить стоп-слова;
- слова в паре объединить символом нижнего подчеркивания “\_”;
- отфильтровать биграммы, которые встретились **не реже** 500 раз (т.е. проводим все необходимые join'ы и считаем NPMI только для них, НО оценку вероятности встретить бигramму считаем на полном датасете);
- отсортировать слова в выводе по значению NPMI;
- вывести **ТОР-39** самых популярных коллокаций и их значения NPMI (округляем до 3го знака после запятой, см. `round`).

Пример вывода.

```
...
south_africa      0.619
roman_empire      0.603
...
```

## 5. Сроки сдачи и правила оформления задания.

**Deadline: 29.07.2019, 08:00**



Оформление задания:

- Код задания (Short name): **HW6:Spark-RDD**.
  - Решения задач должны содержаться в одной папке.
  - Выполненное ДЗ запакуйте в архив **MF2019Q2\_<фамилия>\_HW#.zip** , например -- MF2019Q2\_Ivanov\_HW5.zip. Например, ваше решение лежит в папке my\_solution\_folder, тогда чтобы на Linux и Mac OS создать архив под названием hw.zip ижать его с помощью zip выполните команду<sup>4</sup>:
    - `zip -r hw.zip my_solution_folder/`
- На Windows 7/8/10: необходимо нажать правую кнопку мыши на директорию my\_solution\_folder/, выбрать в открывшемся меню "Отправить >", затем "Сжатая ZIP-папка". Теперь можно переименовать архив.
- Присылайте выполненное задание на почту [bigdata\\_mf2019q2@bigdatateam.org](mailto:bigdata_mf2019q2@bigdatateam.org) с темой письма "Short name. ФИО.". Например: "HW6:Spark-RDD. Иванов Иван Иванович."
  - Перед отправкой письма, оставьте, пожалуйста, отзыв о домашнем задании по ссылке: [http://rebrand.ly/mf2019q2\\_feedback\\_hw05](http://rebrand.ly/mf2019q2_feedback_hw05). Это позволит нам скорректировать учебную нагрузку по следующим заданиям (в зависимости от того, сколько часов уходит на решение ДЗ), а также ответить на интересующие вопросы.

Любые вопросы / комментарии / предложения можно писать:

- в телеграм-канале: [http://rebrand.ly/mf2019q2\\_telegram\\_join](http://rebrand.ly/mf2019q2_telegram_join)
- На почту: [bigdata\\_mf2019q2@bigdatateam.org](mailto:bigdata_mf2019q2@bigdatateam.org)

Всем удачи!

## 7. Дорешка.

Решения после получения фидбека на решение ДЗ можно улучшить. Разрешается одна досылка в течение 1й недели после окончания дедлайна за ДЗ. Соответственно, фидбек за дорешанные ДЗ вы получите в течение 24 часов после окончания deadline следующего ДЗ.

Дорешивать неработающие задания - нельзя. Это позволит исправить дисбаланс между  
    присланными **работающими** заданиями **после** deadline  
VS  
    присланными **НЕработающими** заданиями **до** deadline

---

<sup>4</sup> Флаг -r значит, что будет совершен рекурсивный обход по структуре директории