

1.3 дубликаты, lower, стемминг, лемматизация

- 1 стандартный поиск дубликатов;
 - 2 дубликаты регистр - lower
 - 3 NLTK -стемминг(корни)
 - 4 лемматизация(станд.форма слова)
- рмystem3

0 из практики

Когда все пропущенные/аномалии/дубликаты исправлены - нужен (повторный) чек дубликатов

Для идентификации дубликатов см совпадения по нескольким столбцам

- тем где нет missin val-ues
- в которых маловероятны совпадения!!!

1 стандартный поиск

1 duplicated().sum()
2 value_counts() - наверху будут дубликаты

.drop_duplicates()

2 дубликаты и регистр

.str.lower()

3 стемминг

-нахождение основы слова
(неявные дубликаты)

vs выбирать корни самому
муха мухаммед будут неотличимы
а стемминг отличит

NLTK

1
russian_stemmer =
SnowballStemmer('russian')
2
russian_stemmer.stem(word)

4 лемматизация

-приведение слова к дефолтной форме:

сущ — им падеж, единственное число;

прил — им падеж, единственное число,

мужской род;

глагол, прич, дееприч — глагол в инфинитиве

несоверш вида.

работает с несуществующими словами

pymystem3:

```
from pymystem3 import Mystem
```

```
m = Mystem()
```

```
lemmas = m.lemmatize(text)
```

value_counts слов:

```
from collections import Counter
```

```
print(Counter(lst))
```

case: каких марок телефонов не хватает на сайте(баланс спроса и предложения)

нужно было пересчитать столбец каунт с учетом дубликатов

Нужно определить, по какой марке телефона не хватает предложений на сайте.

Помимо цифр, показывающих, как мало представлено магазинов продающих нехватящую марку,

от вас ждут доказательств, что пользователям действительно интересны малопредставленные магазинами модели.

Изучите поисковые запросы и отзывы, чтобы обосновать свои рекомендации.