

Лабораторна робота № 6

Тема: Бібліотека NLTK мови Python. Використання нормалізованих даних для аналізу тексту. Векторне представлення.

Мета: Відпрацювати практичні навички написання програм на мові Python для обробки текстової інформації використовуючи методи бібліотеки NLTK.

Література: <http://nlp.net/archives/57>
<https://python-scripts.com/matplotlib#5>

Зміст роботи:

Завдання 1.

За матеріалами лекції № 8 виконати завдання:

1. Задати текст з 4 речень виконуючи умови: речення повинні мати різну кількість слів; кілька слів мають повторюватись в різних реченнях.
2. Використовуючи заданий текст створити список унікальних слів («мішок слів»)
3. Розрахувати частоту кожного терміну в документах.
4. Розрахувати зворотню частотність документів.
5. Отримати статистичну міру оцінки важливості слів в документах.
6. Результати роботи оформити в таблицях у текстовому файлі (.docx). Розрахунки доцільно проводити в MS Excel.

Завдання 2.

Використовуючи результати завдання 1 (знайдений мішок слів) написати програму розрахунку **Term Frequency(TF)** - частоти термінів. Результати вивести на екран та у csv –file.

Завдання 3.

Використовуючи результати завдання 1 додати в програму розрахунок **Inverse Document Frequency (IDF)** – зворотньої частотності документа, що вимірює важливість терміна в конкретній колекції документів. Результати вивести на екран та у csv –file.

Завдання 4.

Використовуючи результати завдання 1 та 2 додати в програму розрахунок спеціальної статистичної міри оцінки важливості слова в документі, що є частиною колекції чи корпусу –

Term Frequency - Inverse Document Frequency (TF-IDF).

Результати вивести на екран та у csv –file.

Завдання 5.

Використовуючи отриману статистичну міру оцінки важливості слів в документі побудувати діаграму за допомогою бібліотеки Matplotlib. Результат вивести на екран та у файл.

Завдання 6.

У звіті описати висновки щодо отриманих даних.

Методичні рекомендації.

Модель мішка слів — це спрощене представлення, яке використовується в обробці природної мови та пошуку інформації (IR). У цій моделі текст (наприклад, речення чи документ) представлено як мішок (мультинабір) своїх слів, нехтуючи граматику та навіть порядком слів, але зберігаючи множинність .

Term Frequency (TF) – частота терміну, яка вимірює наскільки часто зустрічається даний термін в обраному документі. Оскільки в великих документах термін буде зустрічатися більшу кількість раз ніж в маленьких, просто кількість знаходжень цього слова нам не вистачає. Тому використовують відносну частоту – відношення числа входження слова до загальної кількості слів в документі.

$$TF = \frac{\text{кількість появ слова у документі}}{\text{загальна кількість слів у документі}}$$

Inverse Document Frequency (IDF) – зворотня частотність документа, що вимірює важливість терміну в конкретній колекції документів. Деякі слова, наприклад прийменники, зустрічаються в усіх документах дуже часто, хоча майже не мають впливу на сенс тексту. Оскільки під час обрахування частоти терміну, ми вважали кожен токен рівнозначним, нам потрібно зменшити оцінку у словах, які присутні у всіх документах. Для цього і обраховують IDF. Його значення відповідає логарифму від відношення загальної кількості документів в колекції, до кількості документів, в яких присутній обраний термін. Такий розрахунок дозволяє додавати ваги словам, які зустрічаються рідко, на противагу тим, що наявні майже у кожному документі.

$$IDF = \log_{10} \frac{\text{кількість документів}}{\text{кількість документів з заданим словом}}$$

Отримані значення **TF** та **IDF** перемножуються для кожного слова і результат використовується у подальшій роботі.

Таким чином ми отримуємо статистичну міру оцінки важливості слова в документі, що є частиною колекції чи корпусу **TF-IDF**.

Контрольні запитання.

1. Поясніть призначення бібліотеки NLTK.
2. Як побудувати список мішка слів?
3. Чи можна для побудови списку мішка слів використати методику токенизації тексту за словами?
4. В чому суть методу Term Frequency(TF)?
5. В чому суть методу Inverse Document Frequency (IDF)?
6. Що є результатом використання частотних методів аналізу текстів?