

Лабораторная работа 8

Тема: Текстовые эмбединги. Работа с предобученными моделями эмбедингов.

Эмбединги – это методы векторизации слов с использованием нейросетей. В отличие от рассмотренных в прошлой работе простых методов векторизации слов, эмбединги рассматривают контекст слова, это позволяет учитывать семантику слов и, таким образом, удобно для поиска семантически близких слов и текстов. Принцип эмбедингов используется во внутренних механизмах современных LLM.

Существует множество моделей текстовых эмбедингов. Ваша задача поработать минимум с двумя моделями и сравнить эффективность их работы на одинаковых задачах.

Сравнение качества эмбеддингов

1. Загрузите предобученные эмбеддинги, например, Word2Vec от Google News, GloVe от Stanford), FastText от Facebook или иные.

Для работы используйте как минимум две различные модели.

2. С использованием библиотеки gensim реализуйте скрипт, который:

- ✓ ищет семантически схожие слова для какого-либо заданного слова;
- ✓ выполняет аналогию типа "Париж – Эйфелева башня, Минск - ?" или "Италия – пицца, Беларусь - ?" или др.
- ✓ **визуализирует векторы 20–30 слов с помощью t-SNE или PCA (**по желанию**).

3. - Сделайте выводы: какая модель лучше справляется с семантическими аналогиями?

Вопросы:

1. Опишите принцип выделения эмбедингов выбранных вами моделей.
2. Какой параметр характеризует семантическую близость слов?
3. Как можно обучить свои текстовые эмбединги и какие библиотеки предоставляют модели для обучения?