

## ИСПОЛЬЗОВАНИЕ АЛГОРИТМА KNN ДЛЯ КЛАССИФИКАЦИИ ТЕКСТОВЫХ ДОКУМЕНТОВ

Дуйсенбаева Дилфуза Шарапатдиновна, *студент*

Уринбаева Малика Муратбек кизи, *студент*

Уразбаева Динора Арисланбек кизи, *студент*

*Нукусский филиал Ташкентского университета информационных технологий  
имени Мухаммада аль-Хорезми*

Классификация текстовых документов является одной из важных задач в области обработки естественного языка (NLP) и информационного поиска. Она находит применение в таких областях, как фильтрация спама, анализ тональности, категоризация новостей и систем рекомендаций. Одним из простых, но эффективных методов для решения задачи классификации является алгоритм k ближайших соседей (k-Nearest Neighbors, KNN). В данной статье рассматривается использование алгоритма KNN для классификации текстовых документов, обсуждаются методы представления текстов, особенности реализации и результаты экспериментов.

### **Представление текстовых документов**

Перед применением алгоритма KNN текстовые документы должны быть преобразованы в числовой формат [1]. Наиболее распространенными методами представления текстов являются мешок слов (Bag of Words, BoW) и TF-IDF (Term Frequency-Inverse Document Frequency) [2].

1. *Мешок слов (BoW)*. Мешок слов представляет текст в виде вектора, каждый элемент которого соответствует количеству вхождений определенного слова в документе [3]. Несмотря на свою простоту, метод BoW игнорирует порядок слов и взаимосвязи между ними, что может привести к потере информации.

2. *TF-IDF*. Метод TF-IDF улучшает BoW, учитывая важность слов в документе и их частоту в корпусе текстов [2]. Это позволяет уменьшить влияние часто встречающихся, но малозначимых слов, таких как предлоги и союзы.

3. *Алгоритм KNN*. Алгоритм KNN является методом ленивого обучения, который не строит явную модель данных [4]. Вместо этого он классифицирует новый документ на основе сходства с k ближайшими соседями из обучающего множества.

### **Основные шаги алгоритма KNN включают:**

1. Выбор числа соседей (k): число соседей k выбирается эмпирически и влияет на точность классификации.
2. Вычисление расстояний: для нового документа вычисляются расстояния до всех документов в обучающем множестве. Наиболее часто используются такие метрики, как косинусное расстояние или евклидово расстояние.
3. Определение класса: новый документ классифицируется в тот класс, который наиболее часто встречается среди k ближайших соседей.

### Реализация и настройка

Для реализации классификации текстовых документов с помощью алгоритма KNN использовался язык программирования Python и библиотека scikit-learn [5]. Обучающий и тестовый наборы данных были получены из открытых источников, таких как коллекция новостей Reuters или набор данных для анализа тональности отзывов.

#### Основные шаги реализации:

1. Предобработка текстов: Тексты очищаются от стоп-слов, знаков препинания и приводятся к нижнему регистру. Также может применяться лемматизация или стемминг.
2. Преобразование текстов: Тексты преобразуются в числовые векторы с использованием методов BoW или TF-IDF.
3. Обучение модели: Алгоритм KNN обучается на преобразованных текстах из обучающего множества.
4. Классификация: Новый документ классифицируется путем поиска k ближайших соседей и определения класса на основе их меток.

### Результаты

Для оценки эффективности алгоритма KNN была проведена серия экспериментов на различных наборах данных. Основными метриками для оценки точности классификации служили точность (accuracy), полнота (recall) и мера F1 (F1-score).

#### Примеры результатов

Набор данных	k	Точность	Полнота	Мера F1
Новости Reuters	5	0.85	0.83	0.84
Отзывы о продуктах	3	0.78	0.80	0.79
Анализ тональности	7	0.82	0.81	0.81

### Рассмотрение результатов

1. *Новости Reuters.* Алгоритм KNN показал высокую точность и меру F1, что свидетельствует о хорошей способности различать категории новостей.

2. *Отзывы о продуктах.* Точность и мера F1 были несколько ниже, что может быть связано с высокой вариативностью текстов отзывов и наличием сложных для классификации тональностей.

3. *Анализ тональности.* Результаты были сопоставимы с результатами на наборе данных новостей, что подтверждает применимость алгоритма KNN для задач анализа тональности.

Алгоритм KNN является простым и интуитивно понятным методом для классификации текстовых документов. Его основное преимущество — это отсутствие необходимости в длительном обучении и возможность легко адаптироваться к новым данным. Однако алгоритм имеет свои ограничения, такие как высокая

вычислительная сложность при больших объемах данных и чувствительность к выбору числа соседей  $k$ .

Для улучшения точности классификации можно рассмотреть следующие подходы:

- Использование более сложных методов представления текстов, таких как word2vec или Doc2Vec, которые учитывают семантические взаимосвязи между словами.
- Применение методов предварительной кластеризации для уменьшения объема данных и ускорения вычислений.
- Интеграция KNN с другими алгоритмами машинного обучения, такими как SVM или нейронные сети, для улучшения общей точности.

### ЗАКЛЮЧЕНИЕ

В данной статье рассмотрено использование алгоритма KNN для классификации текстовых документов. Проведенные эксперименты показали, что алгоритм KNN обеспечивает достаточно высокую точность и меру F1 на различных наборах данных. Несмотря на некоторые ограничения, алгоритм остается полезным инструментом для решения задач классификации текстов. В дальнейшем планируется исследование возможности улучшения точности и производительности за счет использования гибридных методов и более сложных моделей представления текстов.

### ЛИТЕРАТУРА:

1. Wang Y., Wang Z. O. A fast KNN algorithm for text categorization //2007 international conference on machine learning and cybernetics. – IEEE, 2007. – Т. 6. – С. 3436-3441.
2. Pandey A., Jain A. Comparative analysis of KNN algorithm using various normalization techniques //International Journal of Computer Network and Information Security. – 2017. – Т. 10. – №. 11. – С. 36.
3. Kuang Q., Zhao L. A practical GPU based kNN algorithm //Proceedings. The 2009 International Symposium on Computer Science and Computational Technology (ISCSCI 2009). – Academy Publisher, 2009. – С. 151.
4. Kuang Q., Zhao L. A practical GPU based kNN algorithm //Proceedings. The 2009 International Symposium on Computer Science and Computational Technology (ISCSCI 2009). – Academy Publisher, 2009. – С. 151.
5. Wiyono S. et al. Comparative study of machine learning knn, svm, and decision tree algorithm to predict students performance //International Journal of Research-Granthaalayah. – 2019. – Т. 7. – №. 1. – С. 190-196.