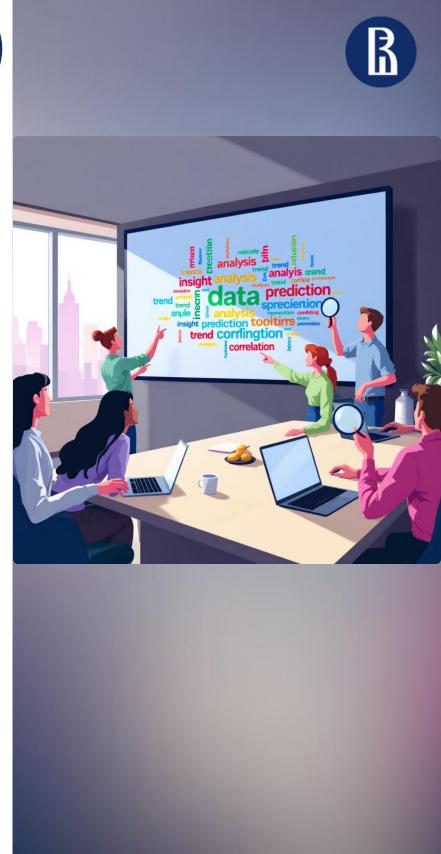
Москва 2025

# Анализ текста (Text Mining)

# Анализ текста (Text Mining) в RapidMiner

Text Mining — технология извлечения и анализа текстовой информации. Используется для поиска закономерностей, классификации документов, выделения ключевых слов и понимания смысловой структуры текстов.



### Значение Text Mining

R

Автоматический анализ текстов используется в маркетинге (отзывы клиентов), финансах (анализ новостей для прогнозов), исследованиях и медицине (выделение симптомов и диагнозов из описаний пациентов).

#### **Text Mining**





### Основные задачи Text Mining

• Классификация и кластеризация текстов.



• Анализ тональности и настроения авторов.

• Извлечение ключевых слов и именованных сущностей для упрощения поиска информации.

### Области применения Text Mining

• Мониторинг социальных сетей и новостных сайтов.

- Анализ обращений и жалоб клиентов.
- Оценка удовлетворенности пользователей продуктами и услугами.









## Этапы работы с текстами

Процесс анализа включает сбор текстов, предобработку (очистка и нормализация данных), преобразование в числовой формат и применение алгоритмов машинного обучения для извлечения выводов.

# Загрузка текстовых данных

Тексты в RapidMiner загружаются с помощью операторов Read Document. Инструмент поддерживает различные форматы, включая txt, docx и PDF, обеспечивая совместимость с разными источниками информации.







# Особенности предобработки текстов





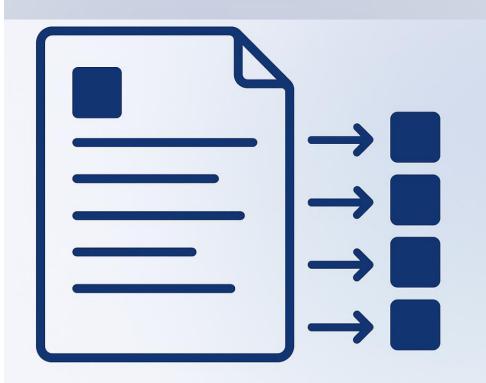


Предобработка — важнейший этап. Включает токенизацию (разделение текста на слова), удаление стоп-слов (часто повторяющиеся и незначимые слова), а также приведение слов к нормальной форме.

#### Токенизация текстов

R

Токенизация — разбиение текста на отдельные элементы (токены). RapidMiner поддерживает разные типы токенизации, например, по пробелам, пунктуации или с использованием регулярных выражений.





#### Стоп-слова и их роль



Стоп-слова (предлоги, союзы, частицы) встречаются часто, но не несут смысловой нагрузки. Их удаление упрощает дальнейший анализ и улучшает точность алгоритмов Text Mining.

#### R

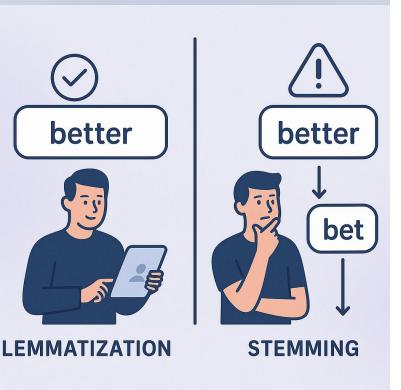
## **Лемматизация и стемминг**

Лемматизация — приведение слов к словарной форме (например, «бегает», «бегают» → «бегать»). Стемминг — отсечение окончаний и суффиксов («бежать», «бегают»  $\rightarrow$  «бега»). Оба метода упрощают обработку текстов.





#### Разница между лемматизацией и стеммингом



Лемматизация учитывает контекст и морфологию языка, более точна, но требует дополнительных ресурсов. Стемминг работает быстрее, но менее точен и может создавать ошибки.

# Генерация векторов признаков

Для анализа текстов в RapidMiner формируются векторы признаков — числовое представление текстов, где каждому слову или термину соответствует определённая числовая частота его появления.









#### Метод TF-IDF

TF-IDF оценивает важность слова относительно конкретного документа и всей коллекции текстов, выявляя значимые термины и игнорируя часто используемые, но малозначимые слова.

#### Применение TF-IDF

R

TF-IDF применяется в информационном поиске, автоматическом аннотировании текстов и создании систем рекомендаций. Позволяет выделить важные слова и темы в больших объёмах информации.





#### Кластеризация текстов



Кластеризация — объединение текстов в группы по смысловому сходству. Алгоритмы, такие как К-Means или иерархическая кластеризация, автоматически выявляют группы схожих текстовых документов.

### B

# Выбор количества кластеров

Оптимальное число кластеров определяется различными методами, например, методом локтя (elbow method) или силуэтом (silhouette analysis), обеспечивая наилучшую интерпретацию полученных групп текстов.





#### Классификация текстов

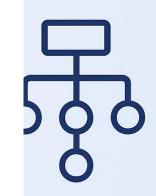


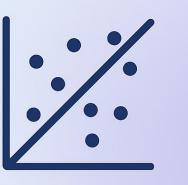
Классификация позволяет автоматически относить тексты к заранее заданным категориям (например, «спам – не спам», «положительный – отрицательный отзыв») на основе обучающих данных.

# Популярные алгоритмы классификации

R

В Text Mining часто применяются алгоритмы: Naive Bayes (быстрый и эффективный), SVM (точный и надёжный), случайный лес (устойчивый к шумам в данных).









# Анализ тональности (Sentiment Analysis)



Анализ тональности оценивает эмоциональный окрас текста (позитивный, негативный, нейтральный). Используется в мониторинге отзывов, социальных сетей и оценке репутации брендов.





### Подходы к Sentiment Analysis

Методы анализа тональности включают использование словарей тональности (эмоционально окрашенные слова) и обучение машинных моделей на размеченных данных (отзывы, рецензии и комментарии).

### Извлечение ключевых слов

Ключевые слова (keywords) отражают основную тему документа. RapidMiner автоматически извлекает такие термины, облегчая последующий анализ и поиск текстовой информации.







### Извлечение именованных сущностей (NER)



NER автоматически определяет имена людей, организаций, географические названия, даты и события в текстах, позволяя быстро выявлять важные факты и структурировать информацию.

## **Автоматизация анализа** текстов

RapidMiner Server позволяет автоматизировать регулярные задачи анализа текстов, выполнять их по расписанию и интегрировать результаты анализа с другими информационными системами.







### Заключение и перспективы Text Mining



Text Mining является важным инструментом анализа больших объёмов текстовой информации. RapidMiner значительно упрощает этот процесс, делая его доступным даже без навыков программирования.