Москва 2025

Web Scraping и анализ данных

Web Scraping и анализ данных в RapidMiner

Texнoлогия web scraping позволяет извлекать и обрабатывать данные с вебстраниц. RapidMiner предлагает инструменты для автоматического сбора, очистки и анализа вебданных без необходимости программирования.



Понятие Web Scraping

Web Scraping — метод автоматического извлечения больших объёмов данных с вебсайтов.

Основная цель – преобразовать неструктурированную вебинформацию в структурированный и пригодный для анализа формат.







Зачем нужен Web Scraping

• Сбор данных о конкурентах и рынке.







• Мониторинг цен на товары и услуги.

• Извлечение данных для исследований и аналитических проектов в различных сферах.

Правовые аспекты Web Scraping

Использование технологии Web Scraping требует соблюдения законодательства о персональных данных, авторском праве и условий использования веб-сайтов. Важно проверять юридическую допустимость извлечения конкретной информации.















Этика и ограничения Web Scraping

- Robots.txt и условий сайта.
- Запрет на массовое копирование защищённых данных.
- Рекомендуется избегать чрезмерной нагрузки на серверы сайтов.

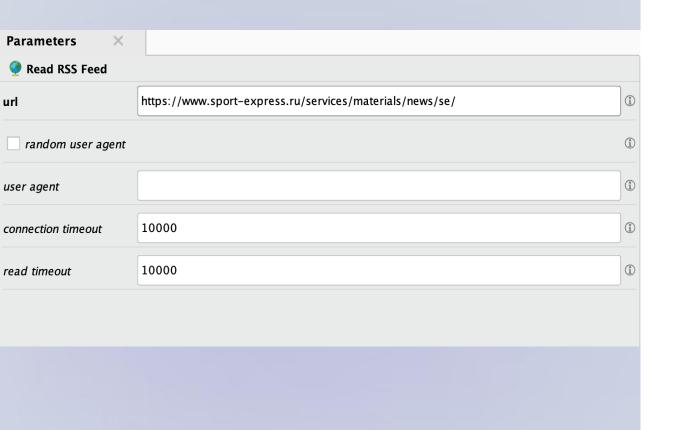
Этапы Web Scraping в RapidMiner

Работа в RapidMiner включает последовательные этапы: подключение к сайту, извлечение информации, очистка данных от ненужных элементов и дальнейший структурированный анализ полученных результатов.









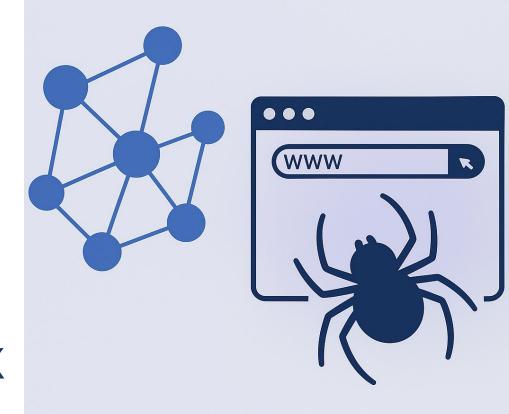
Подключение к вебресурсам

RapidMiner использует операторы Web Mining (Web Crawl, Get Page). Эти инструменты позволяют загружать HTML-страницы и получать контент для дальнейшей обработки и анализа.

Особенности оператора Web Crawl

Оператор Web Crawl предназначен для автоматического обхода множества страниц сайта по ссылкам. Удобен для извлечения данных с больших сайтов, имеющих многоуровневую структуру.







		44					
Row No.	Id	Published	Author	Title	Content	Link	Categories
1	1	Mar 13, 20		Экс-гандбо	Российская	https://www	Гандбол
2	2	Mar 13, 20		Захарян и	Полузащит	https://www	Футбол – Ли
3	3	Mar 13, 20		Галлахер ус	Полузащит	https://www	Футбол – Ли
4	4	Mar 13, 20		13 марта: к	Рассказыва	https://www	Стиль жизн
5	5	Mar 13, 20		Капитан «Б	Капитан «Б	https://www	Футбол – Ли
6	6	Mar 12, 20		Доменикал	Президент	https://www	Авто-мото
7	7	Mar 12, 20		Фанаты «Ат	Фанаты «Ат	https://www	Футбол – Ли
8	8	Mar 12, 20		ЦСКА — «Д	12 марта Ц	https://www	Футбол – Ку
9	9	Mar 12, 20		Файзуллаев	Полузащит	https://www	Футбол – Ку
10	10	Mar 12, 20		Шелтон вы	Американс	https://www	Теннис - АТР
11	11	Mar 12, 20		Полузащит	Полузащит	https://www	Футбол – Ку
12	12	Mar 12, 20		Гасперини	Главный тр	https://www	Футбол – Ит
13	13	Mar 12, 20		Полузащит	Полузащит	https://www	Футбол – Ку
14	14	Mar 12, 20		Койта: «Для	Нападающ	https://www	Футбол – Ку
15	15	Mar 12, 20		Медведев	Российский	https://www	Теннис – АТР
16	16	Mar 12, 20		Жамнов —	Главный тр	https://www	Хоккей – КХЛ
17	17	Mar 12, 20		Mash: проп	Тесть бывш	https://www	Общество
18	18	Mar 12, 20		«Рейнджер	«Рейнджер	https://www	Хоккей – НХЛ

Извлечение элементов HTML

Оператор Extract Information использует XPath-запросы. XPath позволяет выбирать конкретные элементы (теги, классы, идентификаторы) из HTML-страницы и преобразовывать их в структурированный формат.

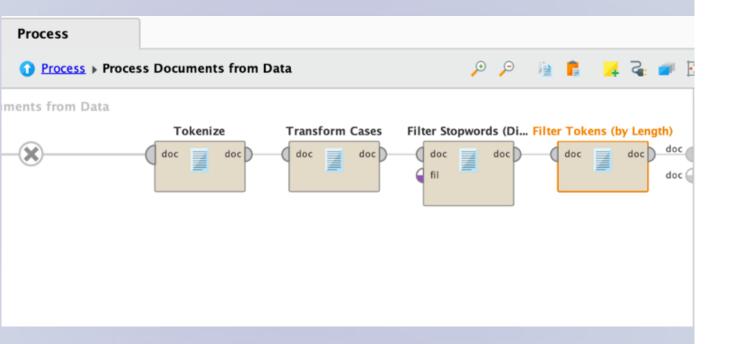
XPath-запросы: краткий обзор

XPath — язык запросов, используемый для навигации по XML и HTML-документам. Например, запрос "//div[@class='price']" извлечёт все элементы div с указанным классом.









Очистка и обработка данных

Полученные данные часто содержат лишние символы, HTML-теги или пробелы. Операторы Replace, Filter Tokens и Trim в RapidMiner помогают быстро привести данные к пригодному для анализа виду.

Типичные проблемы извлечённых данных

• Лишние пробелы и знаки препинания.

• HTML-теги внутри текста.

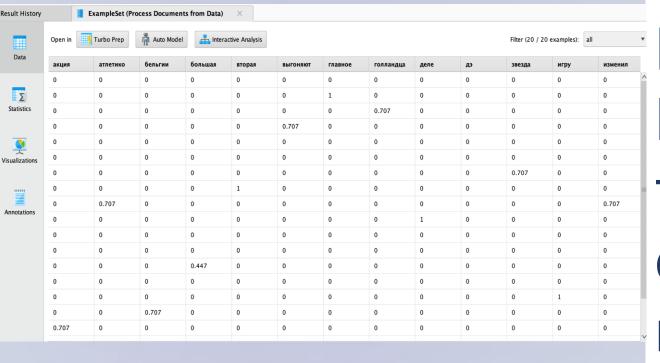
• Некорректные типы данных (текст вместо чисел).







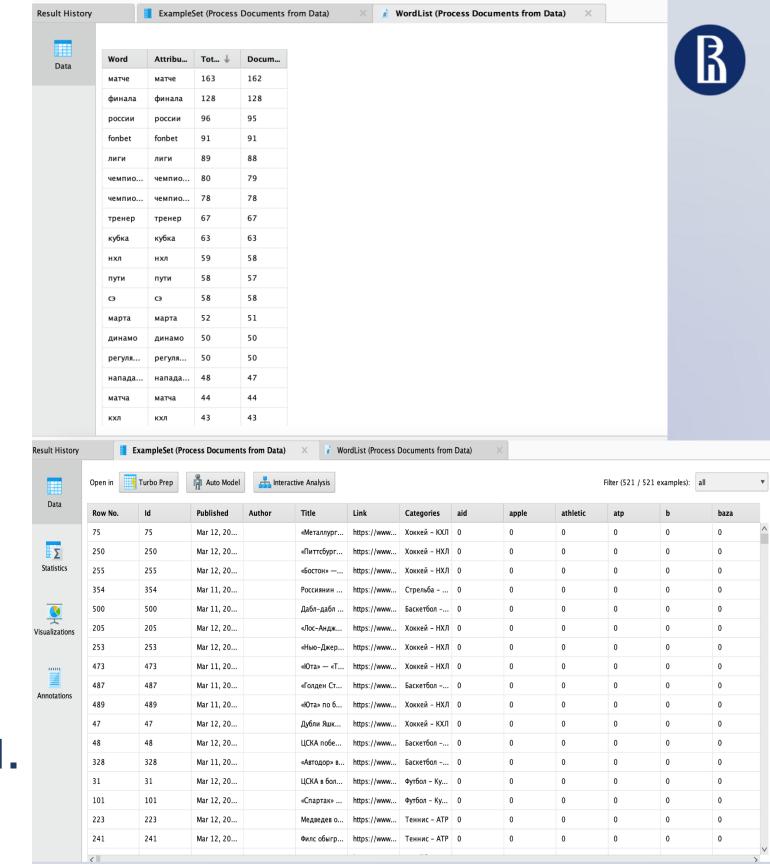
Преобразование типов данных



Использование оператора Parse Numbers преобразует текстовые данные в числовой формат, необходимый для проведения вычислений, фильтрации и последующего анализа.

Анализ данных после сбора

Извлечённые веб-данные подлежат дальнейшему анализу с помощью стандартных инструментов RapidMiner: статистических методов, кластеризации, классификации и прогнозного моделирования.





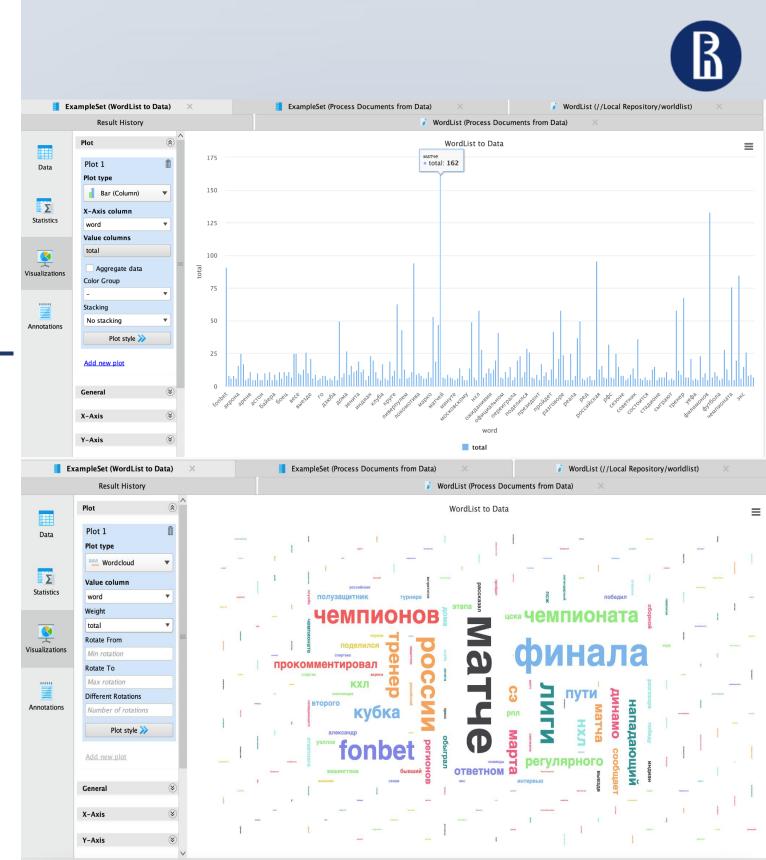
Построение визуализаций

RapidMiner предоставляет инструменты для визуализации собранных данных. Гистограммы, диаграммы рассеивания и столбчатые графики позволяют быстро оценить распределение и структуру информации.

Bar Chart и Word Cloud

На вкладке Charts можно настроить Bar Chart (слово – ось X, частота – ось Y). Самые популярные слова будут на вершине столбцов.

Word Cloud — размер шрифта слова зависит от его частоты.

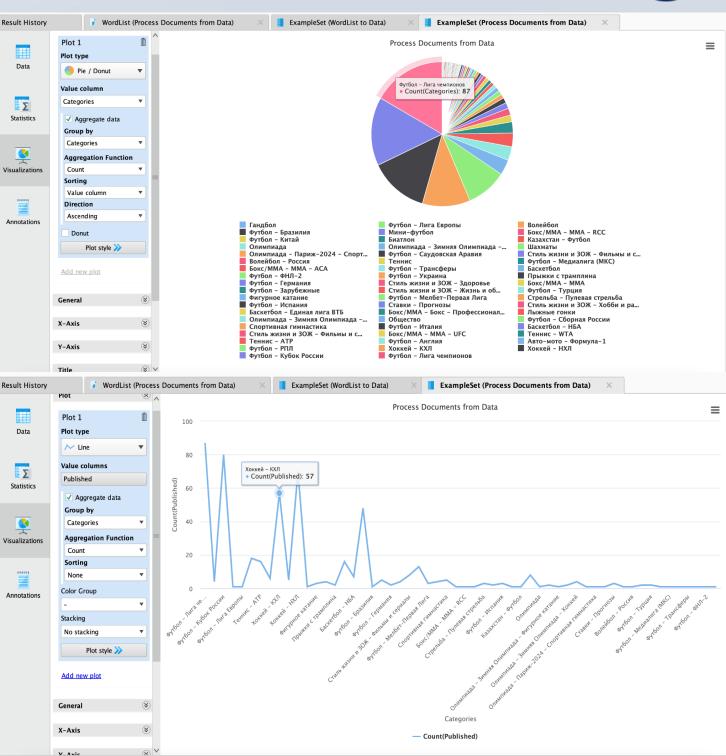


Pie Chart и Line Chart

Если при парсинге удалось получить категорию каждой новости, то можно проанализировать, какие направления спорта преобладают.

Pie Chart помогает понять процентное соотношение. Line Chart подойдет для просмотра динамики во времени.





Диаграммы рассеяния

Scatter/Bubble можно использовать, когда требуется сравнить несколько категорий и упоминаемость набора биграмм.

Подобная диаграмма иллюстрирует в каких рубриках чаще всего всплывают конкретные фразы.



Интерпретация визуализаций

Графики помогают выявить тенденции, аномалии или группировки в собранных данных, облегчая принятие решений и ускоряя процесс анализа больших массивов информации.









Особенности анализа текстовых данных

Текстовые данные требуют специальной обработки: токенизация, удаление стопслов, подсчёт частот слов, лемматизация. Это улучшает точность анализа и интерпретации.

Интеграция RapidMiner с другими сервисами

RapidMiner поддерживает интеграцию с внешними АРІ (Google, Twitter, OpenAI), что позволяет расширять возможности анализа и обработки данных с вебресурсов в режиме реального времени.









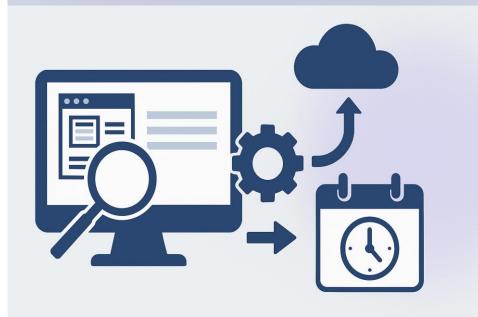
Хранениеизвлечённых данных

Данные после сбора сохраняются в табличных форматах CSV, Excel или базы данных. Это упрощает дальнейшее использование и позволяет интегрировать результаты анализа с другими программами.

Автоматизация процессов Web Scraping

RapidMiner Server позволяет автоматизировать регулярные процессы веб-скрейпинга и анализа данных, выполняя задачи по заданному расписанию, что экономит ресурсы и время.









Проблемы и ограничения Web Scraping

• Частые изменения структуры сайтов.

• Использование JavaScript, затрудняющее извлечение.

• Блокировка IP-адресов при частых запросах к серверу.

Решение сложностей Web Scraping

Использование прокси-серверов, сервисов для обработки JavaScript (например, Selenium) и регулярный мониторинг изменений на сайте помогают эффективно справляться с трудностями веб-скрейпинга.







Примеры успешного применения Web Scraping



Web Scraping используется в маркетинговых исследованиях, мониторинге отзывов клиентов, анализе цен конкурентов и прогнозировании трендов на основе данных из социальных сетей.

Заключение

R

RapidMiner — мощный инструмент для Web Scraping и анализа данных, не требующий написания кода. Грамотное применение методов извлечения и анализа вебинформации значительно повышает эффективность аналитических проектов.

