



6. РАСПОЗНАВАНИЕ ТЕКСТОВ ОБЪЯВЛЕНИЙ И СТРУКТУРИ- РОВАНИЕ РЫНОЧНОЙ ИНФОРМАЦИИ



- Одной из наиболее сложных задач, с которой сталкивается специалист в оценочной деятельности, является структурирование информации, получаемой из текста объявления. Дело в том, что наиболее значимая информация о стоимости объекта оценки заключена в содержании объявления, которое обычно составляется в произвольной форме в соответствии с личными предпочтениями автора объявления (продавца, риелтора и т.п.).
- Для того, чтобы иметь возможность обрабатывать огромные массивы такой информации в автоматическом режиме, ее следует представить в виде таблиц. При грамотной обработке описание объекта, содержащееся в тексте объявления, можно использовать для формирования нужных признаков и представления их **в виде таблиц**, на основе которых выполняется последующая обработка.
- В процессе выполнения данного исследования выделялись различные релевантные признаки. К релевантным данному исследованию отнесены признаки, которые в наибольшей степени влияют на стоимость и ликвидность. Это позволило **представить описание объекта в структурированном виде** и обеспечить дальнейшую обработку этой информации.
- В таблицах ниже представлен перечень ценообразующих факторов, а также значений, которые они принимают, которые в совокупности представляют собой основу базы данных, которая подвергается дальнейшей обработке методами машинного обучения.

Офисная недвижимость

Группа факторов	Ценообразующий фактор		Значения ценообразующих факторов
Общие признаки	Ссылка		
	Источник		
	Телефон		
	Ссылка на скриншот		
	Заголовок объявления		
	Описание		
	Сегмент		Офисная недвижимость
	Тип сделки		Продажа Аренда
Признаки местоположения	Город		
	Точный адрес		
	Координаты		
	Название административного района		
	Территориальная зона (маленькая)		
	Кластер (географический)		
	Кластер (функционально-географический)		
	Территориальная зона (большая)		Исторический центр города
			Центры деловой активности
			Многоквартирная жилая застройка
			ИЖС
			Промзоны
			Окраины
			Зона автомагистралей
			Зеленая зона
			Многофункциональные зоны
			Объекты соц. назначения
	Красная линия		Зона путей морского сообщения
			Прибрежно-курортная зона
			Красная линия
Внутриквартально			
Точки интереса	Ближайшая точка (наименование)	Школа	
		Детский сад	
	Расстояние до точки, напрямую в метрах	Станция метро	
		Остановка общественного транспорта	
	Расстояние до точки, пешком в метрах	Поликлиника/Больница	
		Центр города (Культурный и политический)	
Расстояние до точки, пешком в мин.	Парки/скверы/лесные массивы		
	Море (для курортных регионов)		
	МКАД (для Московской области)		
		КАД (для Ленинградской области)	
		Кладбище	

Офисная недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Тип объекта	Тип объекта	Комплекс зданий	
		Отдельностоящее здание	
Технические характеристики здания	Общая площадь офисного здания, кв.м	Встроенное помещение	
	Общая площадь складского комплекса, кв.м		
	Год постройки		
	Класс		A
			B (B+,B-)
			A+B (если объект высококлассный, но мы не можем отнести ни к A, ни к B)
			C
			объекты вне классификации
	Тип здания		Офисный/Бизнес центр
			Административное здание
			Офисно-торговый комплекс
			Офисно-складской комплекс
			Офисно-жилой комплекс
			Многофункциональный комплекс
			Особняк
	Жилой дом		
	Бизнес центр (наименование)		
	Этажность		
	Материал стен		Кирпичные
			Железобетонные (панельные, блочные)
Монолитные			
Сэндвич-панели			
Доступ		Прочие	
		Свободный	
Охрана		Закрытый (пропускная система)	
		Круглосуточная	
Парковка		Видеонаблюдение	
		Нет	
		Подземная	
		Наземная крытая	
		Наземная открытая (организованная, для легкового транспорта)	
	Стихийная		
	Для грузового транспорта		
	Отсутствует		

Офисная недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Технические признаки помещения	Общая площадь, кв.м		
	Назначение	Офис	
		Прочее	
	Этаж	Первый	
		2 и выше	
		Цоколь	
		Подвал	
		Мансарда	
	Состояние отделки	Первичный рынок	Без отделки
			Под чистовую отделку
			Муниципальный ремонт
			Современный ремонт
		Вторичный рынок	Дизайнерский ремонт
			Требуется капитальный ремонт
			Требуется косметический ремонт
Эконом			
	Современный ремонт		
	Дизайнерский ремонт		
Высота потолков, м			
Планировка	Open-space		
	Смешанная (open-space и коридорная)		
	Коридорная		
Отдельный вход	Есть		
	Нет		

Офисная недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Ценовые признаки	Цена предложения, руб.		
	Удельная цена, руб./кв.м (первоначальная)		
	Удельная цена, руб./кв.м (последняя)		
	Эксплуатационные расходы		
	Коммунальные в цене (для аренды)	Не включено	
		Включено частично	
		Все включено	
	Залог, руб.		
Налоговые платежи			
Временные признаки	Объявление участвует в ликвидности	True False	
	Дата создания		
	Дата парсинга		
	Статус объявления	'white' - когда у объявления отсутствует отметка о том что оно снято или продано на сайте	
		'gray' - когда у объявления присутствует отметка о том что оно снято или продано на сайте и дата последнего парсинга менее 28 дней (4 недель)	
		'black' - когда у объявления присутствует отметка о том что оно снято или продано на сайте и дата последнего парсинга более 28 дней (4 недель)	
	Дата снятия		
	Количество просмотров в неделю		

Торговая недвижимость

Группа факторов	Ценообразующий фактор		Значения ценообразующих факторов
Общие признаки	Ссылка		
	Источник		
	Телефон		
	Ссылка на скриншот		
	Заголовок объявления		
	Описание		
	Сегмент		Торговая недвижимость
	Тип сделки		Продажа Аренда
Признаки местоположения	Город		
	Точный адрес		
	Координаты		
	Название административного района		
	Территориальная зона (маленькая)		
	Кластер (географический)		
	Кластер (функционально-географический)		
	Территориальная зона (большая)		Исторический центр города
			Центры деловой активности
			Многоквартирная жилая застройка
			ИЖС
			Промзоны
			Окраины
			Зона автомагистралей
			Зеленая зона
	Красная линия		Многофункциональные зоны
			Объекты соц. назначения
			Зона путей морского сообщения
			Прибрежно-курортная зона
	Точки интереса		Красная линия
Внутриквартально			
Точки интереса		Школа	
		Ближайшая точка (наименование)	
		Расстояние до точки, напрямую в метрах	
		Расстояние до точки, пешком в метрах	
Точки интереса		Парки/скверы/лесные массивы	
		Море (для курортных регионов)	
		МКАД (для Московской области)	
Расстояние до точки, пешком в мин.		КАД (для Ленинградской области)	
		Кладбище	

Торговая недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Тип объекта	Тип объекта	Комплекс зданий	
		Отдельностоящее здание	
Технические характеристики здания	Общая площадь торгового комплекса, кв.м	Встроенное помещение	
	Общая площадь офисного здания, кв.м		
	Торговая площадь (GLA), кв.м		
	Год постройки		
	Класс		торгово-развлекательные центры и торговые центры;
			районный/микрорайонный торговый центр; стрит-ритейл
			торговые площади
	Тип здания		Многофункциональный комплекс
			Офисно-торговый комплекс
			Торговый комплекс
			Жилой дом
	Наименование ТЦ		Отдельностоящее здание магазина
	Этажность		
	Якорные арендаторы		
	Материал стен		Кирпичные
Железобетонные (панельные, блочные)			
Монолитные			
Сэндвич-панели			
Охрана		Прочие	
		Круглосуточная	
		Видеонаблюдение	
Парковка		Нет	
		Подземная	
		Наземная крытая	
		Наземная открытая (организованная, для легкового транспорта)	
		Стихийная	
		Для грузового транспорта	
Отсутствует			

Торговая недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Технические признаки помещения	Общая площадь, кв.м		
	Назначение		
	Этаж		Первый
			2 и выше
			Цоколь
			Подвал
			Мансарда
	Состояние отделки	Первичный рынок	Без отделки
			Под чистовую отделку
			Муниципальный ремонт
			Современный ремонт
		Вторичный рынок	Дизайнерский ремонт
			Требуется капитальный ремонт
			Требуется косметический ремонт
			Эконом
		Современный ремонт	
		Дизайнерский ремонт	
Высота потолков, м			
Отдельный вход		Есть Нет	
Витринные окна		Есть	
		Нет	

Торговая недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов
Ценовые признаки	Цена предложения, руб.	
	Удельная цена, руб./кв.м (первоначальная)	
	Удельная цена, руб./кв.м (последняя)	
	Эксплуатационные расходы	
	Коммунальные в цене (для аренды)	Не включено Включено частично Все включено
	Залог, руб.	
	Налоговые платежи	
Временные признаки	Объявление участвует в ликвидности	True False
	Дата создания	
	Дата парсинга	
		'white' - когда у объявления отсутствует отметка о том что оно снято или продано на сайте 'gray' - когда у объявления присутствует отметка о том что оно снято или продано на сайте и дата последнего парсинга менее 28 дней (4 недель) 'black' - когда у объявления присутствует отметка о том что оно снято или продано на сайте и дата последнего парсинга более 28 дней (4 недель)
	Статус объявления	
	Дата снятия	
	Количество просмотров в неделю	

Производственно-складская недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Общие признаки	Ссылка		
	Источник		
	Телефон		
	Ссылка на скриншот		
	Заголовок объявления		
	Описание		
	Сегмент	Производственно-складская недвижимость	
	Тип сделки	Продажа Аренда	
Признаки местоположения	Город		
	Точный адрес		
	Координаты		
	Название административного района		
	Территориальная зона (маленькая)		
	Кластер (географический)		
	Кластер (функционально-географический)		
	Территориальная зона (большая)		Исторический центр города
			Центры деловой активности
			Многоквартирная жилая застройка
			ИЖС
			Промзоны
			Окраины
			Зона автомагистралей
			Зеленая зона
			Многофункциональные зоны
			Объекты соц. назначения
Красная линия		Зона путей морского сообщения	
		Прибрежно-курортная зона	
Точки интереса		В непосредственной близости от автомагистрали	
		На удалении от магистрали	
		Автомагистраль	
		Железнодорожные пути	
Точки интереса		Центр города	
		МКАД (для Московской области)	
		КАД (для Ленинградской области)	

Производственно-складская недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов
Тип объекта	Тип объекта	Комплекс зданий
		Отдельностоящее здание
Технические характеристики здания	Общая площадь здания, кв.м	Встроенное помещение
	Общая площадь складского комплекса, кв.м	
	Размер офисной части, кв.м	
	Размер прилегающей территории, га. (Площадь земельного участка, га)	
	Год постройки	
	Физическое состояние здания	Хорошее
		Удовлетворительное
		Неудовлетворительное
	Класс	A
		B (B+, B-)
		A+B (если объект высококлассный, но мы не можем отнести ни к A, ни к B)
		C
		объекты вне классификации
	Тип здания	Складской комплекс
		Производственно-складской комплекс
		Офисно-складской комплекс
		Производственный комплекс
		Склад
		Производство
Многофункциональный комплекс		
Отдельностоящее здание		
Этажность	Ангар	
Материал стен	Кирпичные	
	Железобетонные (панельные, блочные)	
	Монолитные	
	Сэндвич-панели	
	Металлические	
Огороженная территория (забор)	Прочие	
	Есть	
Наличие Ж/Д ветки (на участке)	Нет	
	Есть	
Доступ	Нет	
	Свободный	
	Закрытый	

Производственно-складская недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов	
Технические характеристики здания	Охрана	Круглосуточная Видеонаблюдение	
	Автоматические ворота	Есть Нет	
Технические признаки помещения	Общая площадь, кв.м		
	Назначение		
	Этаж		Первый 2 и выше Цоколь Подвал Мансарда
		Состояние отделки	Хорошая Удовлетворительная Неудовлетворительная
			Высота потолков, м
		Материал пола	
	Материал перекрытий		Бетон Кирпич Смешанный
			Отопление
	Мощность электросети, кВт		
	Водоснабжение		Есть Нет
		Количество мокрых точек, шт.	
	Грузоподъемной механизм		Есть Нет
		Холодильная камера	

Производственно-складская недвижимость

Группа факторов	Ценообразующий фактор	Значения ценообразующих факторов
Ценовые признаки	Цена предложения, руб.	
	Удельная цена, руб./кв.м (первоначальная)	
	Удельная цена, руб./кв.м (последняя)	
	Эксплуатационные расходы	
	Коммунальные в цене (для аренды)	Не включено
		Включено частично
		Все включено
Залог, руб.		
Налоговые платежи		
Временные признаки	Объявление участвует в ликвидности	True False
	Дата создания	
	Дата парсинга	
	Статус объявления	'white' - когда у объявления отсутствует отметка о том что оно снято или продано на сайте
		'gray' - когда у объявления присутствует отметка о том что оно снято или продано на сайте и дата последнего парсинга менее 28 дней (4 недели)
		'black' - когда у объявления присутствует отметка о том что оно снято или продано на сайте и дата последнего парсинга более 28 дней (4 недели)
	Дата снятия	
Количество просмотров в неделю		

АЛГОРИТМ ОБРАБОТКИ НЕЧИСЛОВЫХ ПАРАМЕТРОВ

Гипотеза распределения. Слова, которые встречаются рядом друг с другом и используются в одних и тех же контекстах, имеют тенденцию иметь схожие значения.

ЭТАПЫ ОБРАБОТКИ ПРИЗНАКА



СБОР ДАННЫХ

Domofond

- link
- description

Move

- link
- opisanie = description

Afy

- link
- description

Avito

- link
- description

Количество объектов в объединенной выборке: **более 2 млн. записей**

ОЧИСТКА ТЕКСТА

- 1) Удаление url-адресов
- 2) Удаление специальных идеограмм, пиктограмм и смайликов
- 3) Транслитерация латинских символов в кириллицу по определенному стандарту (ГОСТ 7.79-2000, ГОСТ Р 52290-2004, ГОСТ Р 7.0.34-2014, ICAO DOC 9303, стандарт Википедии основанный на BGN/PCGN). Также есть возможность замены символов, у которых очень схожее написание. Например, латинская буква "с" и "с", написанная на кириллице
- 4) Замена знаков восклицания и вопросительного на знак "."
- 5) Замена схожих по написанию знаков тире на единый формат "-"
- 6) Нормализация нестандартных символов
- 7) Функция для замены несколько подряд идущих символов пробела, знаков табуляции на один знак пробела
- 8) Преобразование чисел формата "сто пятьдесят два" на числовое значение "152"
- 9) Приведение текстов объявлений к нижнему регистру
- 10) Функция для замены орфографических ошибок и опечаток слов с дополнительными преобразованиями для анализа

ВЫДЕЛЕНИЕ ОСНОВНОГО ИНФОРМАТИВНОГО ПРЕДЛОЖЕНИЯ И ПРИВЕДЕНИЕ В НАЧАЛЬНУЮ ФОРМУ

Функция выделения основного предложения работает на поиске опорных слов в предложении:

```
"(?:\S*ремонт\S*|\S*отделк\S*|\S*реконструкц\S*)|(?:\S*Ремонт\S*|\S*Отделк\S*|\S*Реконструкц\S*)"
```

	description	link	sentences	true_sentence
6	Аренда . Офисные кабинеты в ЖК Олимп . Распол...	https://bashkortostan.afy.ru/ufa/snyat-ofis/60...	[Аренда, Офисные кабинеты ЖК Олимп, Расположен...	Расположены на втором этаже Общи вход этажа П...
10	Описание Бизнесцентра Удобные выезды на Дмитр...	https://afy.ru/moskva/snyat-ofis/60001208784	[Описание Бизнесцентра Удобные выезды на Дмитр...	Отделка офиса предусмотрена индивидуально по з...

Приведение в начальную форму:

- 1) Проверка и исправление правописания в предложении
- 2) Удаление повторяющихся слов
- 3) Приведения слов в начальную форму, исходя из текста: Отделка офиса -> отделка офис
- 4) Удаление стоп-слов, кроме несущих смысловую нагрузку: частица не, нет и др.

	description	link	sentences	true_sentence	initial_sentence
6	Аренда . Офисные кабинеты в ЖК Олимп . Распол...	https://bashkortostan.afy.ru/ufa/snyat-ofis/60...	[Аренда, Офисные кабинеты ЖК Олимп, Расположен...	Расположены на втором этаже Общи вход этажа П...	расположить второе этаж общий вход этаж план п...
10	Описание Бизнесцентра Удобные выезды на Дмитр...	https://afy.ru/moskva/snyat-ofis/60001208784	[Описание Бизнесцентра Удобные выезды на Дмитр...	Отделка офиса предусмотрена индивидуально по з...	отделка офис предусмотреть индивидуально зака...

ВЫДЕЛЕНИЕ ОСНОВНОГО ИНФОРМАТИВНОГО ПРЕДЛОЖЕНИЯ И ПРИВЕДЕНИЕ В НАЧАЛЬНУЮ ФОРМУ

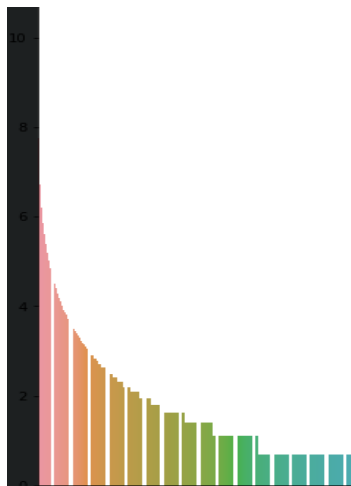
Пример преобразования начального варианта текста (слева) в конечный формат (справа)

<p>ID:9791 Предлагаются в аренду отапливаемые склады «В» класса. Общей площадью 3 208 м² + 2 777 м² в разных помещениях на 2 этаже. Рабочая высота 4,3 м. Ворота докового типа - 6, грузовые лифты г/п 3-5 тонн. Полы бетон с <u>антипылевым</u> покрытием. Режим работы 24/7. Охраняемая парковка, телефония и интернет, система видеонаблюдения. Эксплуатационные платежи включены в ставку.</p>	<p>ид : 9791 предлагаются в аренду отапливаемые склады « в » класса . общей площадью 3208 кв . м . плюс 2777 кв . м . в разных помещениях на 2 этаже . рабочая высота 4 , 3 м . ворота докового типа - 6 , грузовые лифты г / п 3-5 тонн . полы бетон с <u>антипылевым</u> покрытием . режим работы 24/7 . охраняемая парковка , телефония и интернет , система видеонаблюдения . эксплуатационные платежи включены в ставку .</p>
<p>ОФИСНЫЕ ПОМЕЩЕНИЯ ОТ СОБСТВЕННИКА! При проектировании офисного комплекса мы <u>тщательно исследовали</u> все детали, которые важны арендаторам и <u>устранили</u> все мелочи, которые могут вызвать <u>малейший дискомфорт</u>. В этом офисе вам будет максимально комфортно. ✓ Большая парковка с ул. Уральской и на территории, свободные места гарантируем ✓ 2 этаж, охрана, видеонаблюдение ✓ Вся необходимая для комфортной работы инфраструктура ✓ Звукоизолированные помещения ✓ Комфортный микроклимат в помещении благодаря вентиляции и кондиционированию ✓ Общая комната отдыха с кухней на этаже ✓ Бесплатный <u>wi-fi</u> в здании отличный LTE-сигнал ✓ Возможность расширения помещения ✓ Рядом с СБС, очень интенсивный <u>трафик</u> Офис состоит из двух сквозных комнат по 22,4 кв.м. №10/2 и 10/3 на <u>плане</u> <u>звоните!</u></p>	<p>офисные помещения от собственника . при проектировании офисного комплекса мы <u>тщательно исследовали</u> все детали , которые важны арендаторам и <u>устранили</u> все мелочи , которые могут вызвать <u>малейший дискомфорт</u> . в этом офисе вам будет максимально комфортно . большая парковка с <u>ул . уральской</u> и на территории , свободные места гарантируем 2 этаж , охрана , видеонаблюдение <u>вся необходимая для комфортной работы инфраструктура</u> <u>звучо изолированные помещения</u> <u>комфортный микроклимат в помещении</u> благодаря <u>вентиляции и кондиционированию</u> <u>общая комната отдыха с кухней на этаже</u> <u>бесплатный вй-фи</u> <u>плюс в здании отличный дте-сигнал</u> <u>возможность расширения помещения</u> <u>рядом с сбс</u> , очень интенсивный <u>трафик</u> <u>офис состоит из 2 сквозных комнат по 22 , 4 кв . м . до 10/2 и 10/3 на плане</u> <u>звоните .</u></p>

Выделение и подсчет n-gram(n = 2)



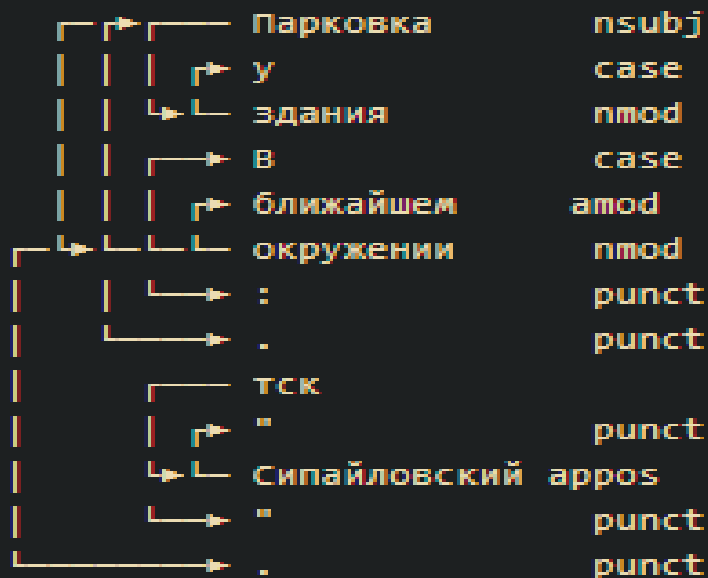
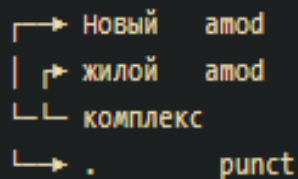
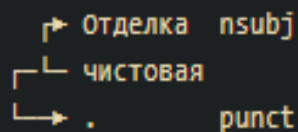
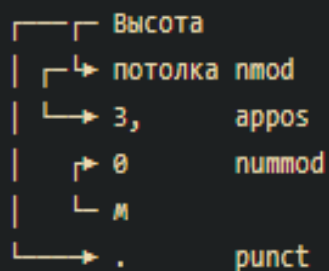
Computational Linguistics



```
(( 'помещение', 'ремонт'), 24122),  
(( 'хороший', 'ремонт'), 23719),  
(( 'косметический', 'ремонт'), 17116),  
(( 'качественный', 'ремонт'), 16257),  
(( 'чистовой', 'отделка'), 12749),  
(( 'помещение', 'отделка'), 11962),  
(( 'офисный', 'отделка'), 9314),  
(( 'свежий', 'ремонт'), 8985),
```

ВЫДЕЛЕНИЕ НУЖНЫХ СТРУКТУР

С помощью синтаксических связей можно получить все связанные слова с опорным словом для дальнейшего подсчета.



ВЫДЕЛЕНИЕ НУЖНЫХ СТРУКТУР

