

Извлечение атрибутов из объявлений на hh.ru

Олег Мельников, 3 курс МКН

**Научный руководитель:
Сергей Резник, Яндекс**

Проблема

- В теории, у абстрактного объявления ожидается наличие строгих характеристик и неформального текста
- Фактически, в ведущих сервисах объявлений о работе строгая часть не очень проработана - в отличие от недвижимости, авто, гостиниц и т.п.
- hh.ru: поиск по строгим атрибутам не очень мощный, поиск по тексту тоже не достаточен (можно искать по слову Scala, но в выдачу попадет и вакансия, где это ключевое требование, и та, где это nice-to-have)
- Для аналитики тенденций на рынке труда хочется вытаскивать полезную информацию из неформальных описаний вакансий

Как это выглядит

- <https://spb.hh.ru/vacancy/38487280>

```
1 <p><strong>Страховой Дом ВСК</strong> осуществляет страховую деятельность с 11 февраля 1992 года Компания входит в ТОП
2 | -10 компаний лидеров на рынке страхового бизнеса. Общий штат более 7000 сотрудников по всей России. В ИТ
3 | подразделении уже более 500 сотрудников.</p>
4 <p>Наши Сотрудники – это единая команда, поэтому взаимодействие внутри коллектива строится на принципах взаимного
5 | уважения и профессионализма.</p>
6 <p>Сейчас мы находимся в поиске <strong>ведущего администратора Linux </strong>сразу на несколько проектов, один из них:
7 </p>
8 <p><strong>Стратегическая программа «Прорыв ДМС» -</strong> это выход на новый уровень сервиса для клиентов компании по
9 | добровольному медицинскому страхованию. Работа в кроссфункциональных командах по agile-принципам. Решение
10 | амбициозных задач по разработке высокотехнологичных инструментов для взаимодействия с клиентами и партнерами.</p>
11 <p><strong>Ключевые задачи:</strong></p>
12 <ul>
13 | <li>Построение непрерывной интеграции и доставки (CI/CD);</li>
14 | <li>Администрирование и настройка Linux-систем (RHEL, CentOS, Debian);</li>
15 | <li>Администрирования HTTP-серверов (Nginx, Apache HTTP Server);</li>
16 | <li>Написание и работа с SQL-запросами (PostgreSQL, MySQL);</li>
17 | <li>Администрирование серверов приложений (PHP-FPM, PHP, Node.js, Redis, RabbitMQ);</li>
18 | <li>Опыт работы с системами контроля версий git;</li>
19 | <li>Мониторинг доступности, производительности и создание проверок для системы мониторинга (Zabbix);</li>
20 | <li>Автоматизация задач на скриптовых языках (Bash, sh);</li>
21 </ul>
22 <p><strong>Для нас важно:</strong></p>
23 <ul>
24 | <li>Знание Linux-систем на уровне администратора (RHEL, CentOS);</li>
25 | <li>Опыт построения непрерывной интеграции и доставки (CI/CD);</li>
```

Задача

По набору текстовых объявлений с hh.ru (для простоты ограничиться IT-сегментом) извлекать:

- Перечень требуемых компетенций
- Перечень желательных компетенций
- Есть ли гибкое расписание (если в тексте сказано с “работа в офисе с 10 до 18”, установить этот атрибут в False)
- Подобласть (Разработка, Тестирование, DevOps, Data Science)

Pipeline - 1

- Скачиваем много объявлений разных IT вакансий (Python, C++, Java, ...)
- Парсим сырое тело объявления на разделы

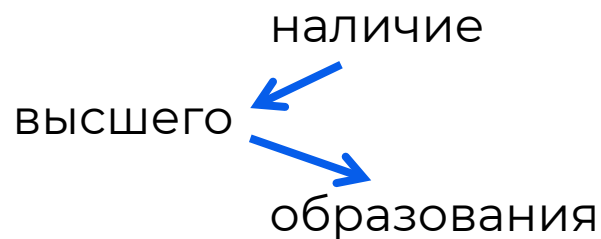
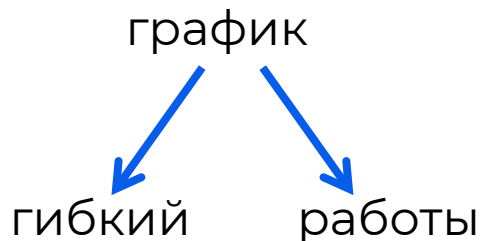
```
32 <p><strong>Плюсом будет:</strong></p>
33 <ul>
34   <li>Знание Agile подходов.</li>
35   <li>Опыт написания скриптов на языке Python;</li>
36   <li>Администрирование и опыт работы с Apache ActiveMQ, Tomcat;</li>
37   <li>Умение читать код.</li>
38 </ul>
```

Pipeline - 2

- Очищаем текст (tokenization, stopwords, lemmatization)
- Обучаем модель распознавания разделов
 1. Чем предстоит заниматься/ключевые задачи
 2. Требования/мы ждем от вас
 3. Nice-to-have
 4. Что мы предлагаем/почему мы/общее описание команды/плюшки от компании
 5. Условия работы
 6. Текущий стек

Pipeline - 3

- Извлекаем требования (NER, POS-tagging)
Python – upos: PROPN, ner: S-MISC
- Извлекаем информацию про график работы (Dependency parsing)



- Извлекаем информацию о подобласти (эвристики)

Результат работы

Вакансия



```
{...} result.json > ...
1  {
2    "Обязательные компетенции": [
3      "Windows",
4      "CI",
5      "CD",
6      "Node.js",
7      "RHEL",
8      "PHP",
9      "Redis",
10     "Linux",
11     "MySQL",
12     "Bash",
13     "Zabbix",
14     "CentOS",
15     "Администрирование серверов Linux",
16     "PostgreSQL",
17     "Nginx"
18   ],
19   "Желательные компетенции": [
20     "Tomcat",
21     "Agile",
22     "ActiveMQ",
23     "Python",
24     "Apache"
25   ],
26   "Гибкий график работы": false,
27   "Подобласть": "DevOps"
28 }
```


Оценка качества

- Качество классификации разделов
CV accuracy = 78%, baseline = 38%
- Качество ключевого алгоритма
методом сэмплирования
Precision – 92%, Recall – 82%

Используемые технологии

- Для выкачки данных и парсинга:



BeautifulSoup

- Для обработки текста:



Ссылки

- Мой репозиторий -
<https://github.com/melnikoff-oleg/hh-load>
- Stanza -
<https://stanfordnlp.github.io/stanza/index.html>

Спасибо за внимание

Контакты

- telegram: @melnikoff_oleg
- mail: oleghools@gmail.com