

# Задача

1). Сделать настройки [программы Datacol](#) для парсинга объявлений с [сайта Avito](#) в таблицу MySQL по заданным требованиям:

- Регион :: **Санкт-Петербург, Ленинградская область, Татарстан**
- Тип объявления :: **только “Частные”**
- Категории/подкатегории для **стартовых URLOв** ::
  - Транспорт > Автомобили
  - Недвижимость > Квартиры
  - Недвижимость > Комнаты
  - Недвижимость > Дома, дачи, коттеджи
  - Недвижимость > Земельные участки
  - Недвижимость > Гаражи и машиноместа
  - Недвижимость > Коммерческая недвижимость
  - Недвижимость > Недвижимость за рубежом
  - Личные вещи > Часы и украшения
  - Личные вещи > Красота и здоровье
  - Транспорт > Мотоциклы и мототехника
- Категории/подкатегории, с которыми работаем **в принципе** :: Все, кроме “Работа > Вакансии” и “Работа > Резюме” и “Хобби и отдых > Знакомства”

2). Спроектировать таблицу MySQL, в которую будут складываться данные. Примерный формат:

- все значащие поля (см. описание ниже);
- текстовое поле для всего объявления;
- дата создания объявления;
- дата обновления объявления.

Номер объявления первичный или уникальный ключ.

## Формат URL запроса на поиск нужных нам объявлений. Проверки.

Соответственно, URLы разделов, к которым мы станем обращаться, будут вида:

```
https://www.avito.ru/tatarstan/kvartiry?user=1  
https://www.avito.ru/sankt-peterburg/avtomobili?user=1  
https://www.avito.ru/leningradskaya_oblast/chasy_i_ukrasheniya?user=1
```

Где:

- **user=1** – это “частные” объявления
- **tatarstan, sankt-peterburg, leningradskaya\_oblast** – это искомые регионы
- **kvartiry** и т.п. – это категории/подкатегории

Допустимо использование мобильной версии сайта для решения нашей задачи, если это будет проще без ущерба качеству итоговых данных.

При написании настроек необходимо учесть, что Datasol очень часто при парсинге “уходит” в другие категории, регионы и типы объявлений. Причем “уйти” он может как на страницу со *списком* объявлений, так и напрямую на страницу *конкретного* объявления (через блок “похожие объявления”).

Для нас *недопустим* уход в иные регионы кроме заданных, а также *критически недопустим* переход от “частных” объявлений к объявлениям компаний.

Поэтому необходимо сделать в настройках следующие проверки:

- проверку на нужный регион - в URLe должны содержаться **tatarstan/sankt-peterburg/leningradskaya\_oblast**
- проверку на нужный тип (“частные”) объявления
  - для страницы со *списком* объявлений - в URLe должен содержаться **user=1**
  - на странице *конкретного* объявления в URLe тип юзера не указан, поэтому проверять надо как-то иначе.

**Насколько я понимаю**, логика обозначения НЕчастного объявления на его странице такая:

- либо в строке `<div class="description_seller"> ... </div>` указан в значении НЕ “Продавец”, а что-то иное - “Автодилер” для автомобилей ([пример](#)), “Агенство” для недвижимости ([пример](#)) и т.п.
- либо в строке `<div class="description_seller"> ... </div>` указан все-таки “Продавец”, но при этом справа от его имени есть бейдж вида `<span class="с-2"> (...)</span>` со значениями типа “компания”, “магазин” и т.п. ([пример](#))
- **возможно, есть и иные признаки НЕчастного объявления, о которых я не упомянул**

## Перечень полей для парсинга. Формат значений.

ПАРАМЕТР(Ы)	КАК СПАРСИТЬ / ПРИМЕР ЗНАЧЕНИЯ
Категория - Подкатегория объявления	Берем из каталога. Примеры: Категория :: Недвижимость Подкатегория :: Комнаты  Категория :: Транспорт Подкатегория :: Автомобили

<b>Тип объявления</b>	<p>Варианты значений:</p> <ul style="list-style-type: none"> <li>Продам</li> <li>Сдам</li> <li>Куплю</li> <li>Сниму</li> </ul> <p>Если не указано - присваивать значение по умолчанию "Продам"</p>
<b>Номер (ID) объявления</b>	<p>Брать как со страницы самого объявления...</p> <p>Номер объявления: 587866915</p> <p><b>Важно! Для случаев, когда Datasol повторно парсит объявление с уже существующим в нашей базе ID:</b> Таблица у нас в MySQL одна. В ней номер объявления является первичным или уникальным ключом. При вставке объявления первый раз проставляется дата вставки. При повторной вставке этого же объявления по номеру, идет запрос <code>insert on duplicate key update</code> и проставляется дата обновления.</p>
<b>URL объявления</b>	<p><code>https://www.avito.ru/moskva/avtomobili/honda_cr-v_2011_587866915</code></p>
<b>Дата размещения объявления на сайте</b>	<p>Можно брать из списка объявлений или со страницы самого объявления.</p> <p>Формат:</p> <ul style="list-style-type: none"> <li>(Размещено) сегодня в ЧЧ:ММ</li> <li>(Размещено) вчера в ЧЧ:ММ</li> <li>(Размещено) ДД месяца в ЧЧ:ММ</li> </ul> <p>Соответственно, для нашей БД конвертим в формат нормальной даты /времени <code>ДД.ММ.ГГГГ ЧЧ:ММ</code></p>
<b>Наименование товара</b>	<p>Берется из <code>title</code> объявления из списка объявлений или со страницы объявления.</p>
<b>Цена (руб.)</b>	<p>Берем из соответствующего поля <code>Цена</code> конкретного объявления</p> <p><b>Важно!</b> Цену нужно доставать одну (их там может быть две) и сразу в виде числа (удаляя пробелы и "руб.")</p>
<b>Имя пользователя</b>	<p>В зависимости от категории, лейбл у поля на странице объявления может различаться ("Меня зовут", "Продавец", ...), но в коде это, похоже, всегда</p> <pre>&lt;div class="description_content" id="seller" itemprop="seller" itemscope="" itemtype="http://schema.org/Person"&gt;   &lt;strong itemprop="name"&gt;     <b>Вазген</b>   &lt;/strong&gt; &lt;/div&gt;</pre>
<b>Телефон пользователя</b>	<p>Берем значение со страницы конкретного объявления - расположен под именем</p>

	<p>Телефон частично скрыт от визуального считывания за кнопкой “Показать телефон” - придется, видимо, эмулировать клик по кнопке:</p> <pre data-bbox="682 224 1885 289">&lt;span class="description__phone-insert js-phone-show__insert"&gt;&lt;span class="btn__text"&gt;Показать телефон&lt;/span&gt;&lt;/span&gt;</pre> <p><b>Важно!</b> Телефон необходимо в базу класть в виде числа (удаляя пробелы и дефисы)</p>
<b>Страна</b>	<p>Так как avito.ru работает только с россиянами, для всех пишем по умолчанию - Россия</p>
<b>Область - Город - Метро</b>	<p>Если смотрим в поле “Город” на странице конкретного объявления, то логика такова:</p> <p>Просто текст = &lt;Город&gt; (областной центр)</p> <p>Текст, Текст = &lt;Область/регион&gt;, &lt;Город&gt; (не центральный) <a href="#">ПРИМЕР</a></p> <p>Текст, <b>м.</b> Текст = &lt;Город&gt; (областной центр), м. &lt;Станция метро&gt; <a href="#">ПРИМЕР</a></p> <p>Иных вариантов я не встречал. Может быть вы сможете парсить более изящным способом, но на выходе нам нужны значения трех полей:</p> <p>Область</p> <p>Город</p> <p>Метро (если есть)</p>
<b>HTML-текст всего объявления</b>	<p>Весь значащий текст объявления.</p>