

Спайдер сбора ссылок комментариев для создания своей тематической базы дропов под названием Black Widow Spider. Сбор своей базы дропов по своей тематике. Существует обмен комментариями сайтов по своей тематике, например, на сайте с рецептами комментируют люди, имеющие близкую тематику и оставляют ссылку, ведущую на сайт, это практически сообщество, довольно узкий круг, без залетных.

Для чего это нужно?

Многие сайты уходят в небытие - люди забрасывают сайты, переезжают, меняют работы - и это, как правило, очень хорошие сайты с авторскими текстами - я работаю именно так, только отбираю вручную, отдаю программе и работаю с этими забытыми доменами, руки не доходят автоматизировать, позже приведу пример, если не ясно. А просто скаченный список доменов по ключевым словам - это 80% мусора, который программа перелопачивает, прежде чем наткнутся на что-то стоящее.

Что он делает?

Паук ходит по тематическим сайтам (например, строительной, женской или сайтам другой тематики) и собирает ссылки, находящиеся в никах пользователей, ведущей на их сайт.

Как работает?

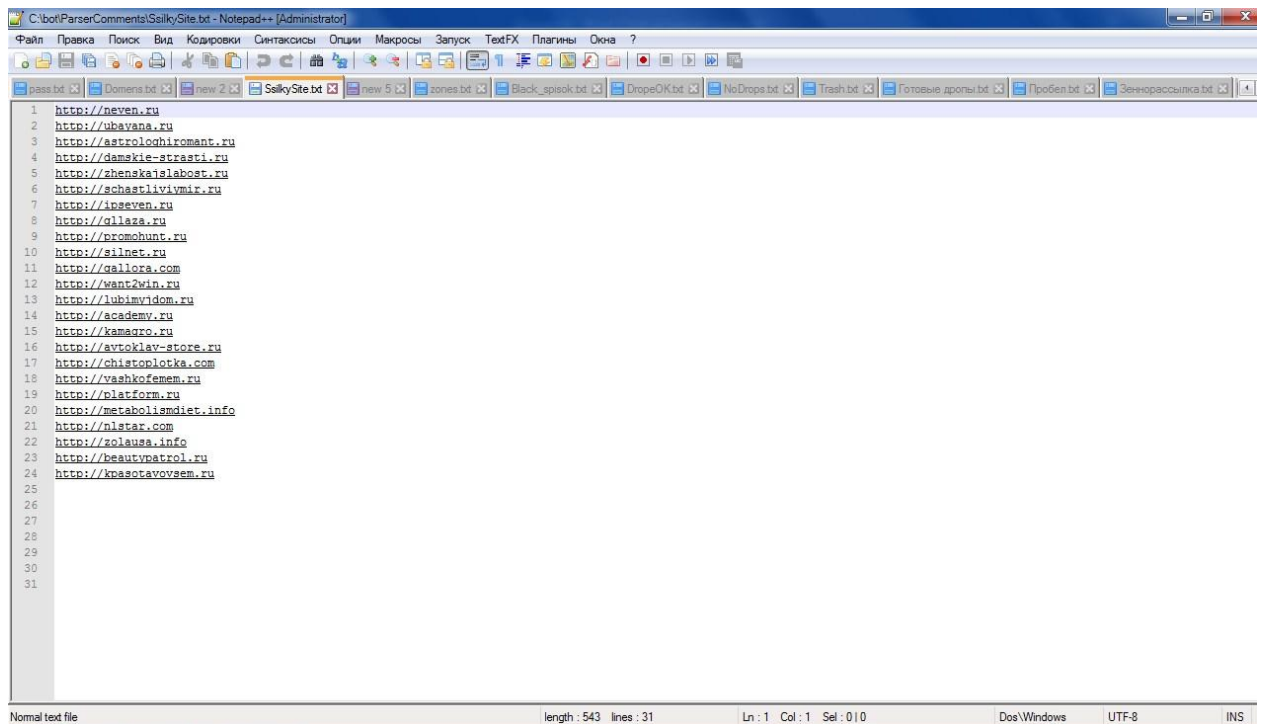
1. БОТ берет ссылки с любых сайтов на любой платформе и на любых языках
2. БОТ самообучается - во время поиска ссылок пишет данные всех сайтов, на которые заходил. Если он зашел на Фейсбук или Твиттер и ему там не понравилось, при следующем случайном заходе он сравнит список и на них не пойдет.
3. Примерный принцип работы БОТА (алгоритм очень сложный) - Бот берет данные из файла (начальные URL, которые вы ему скормите, смотрите видео), создает карту сайта из данного урл и начинает ходить по страницам в поисках ссылок, берет ссылки из комментариев и начинает их проверять на ответ сервера, если ответа нет, бот идет в Вебархив и проверяет, есть ли там сайт. Если ответ положительный, смотрит количество ссылок, и если считает, что сайт перспективный, заносит его в список дропов.

Эта программа идет в связке с WebArchiveMasters - спайдер собирает тематические дропы, а парсер берет с них текста. Теперь не нужны списки доменов, программа найдет всё сама.

Что нужно сделать для работы паука.

Ввести в поисковый ключевой запрос по вашей тематике - так-как я работаю по женской тематике, я ввожу запрос: "женский сайт свежие комментарии". После этого пройтись по выдаче и найти 5-6 сайтов с комментариями.

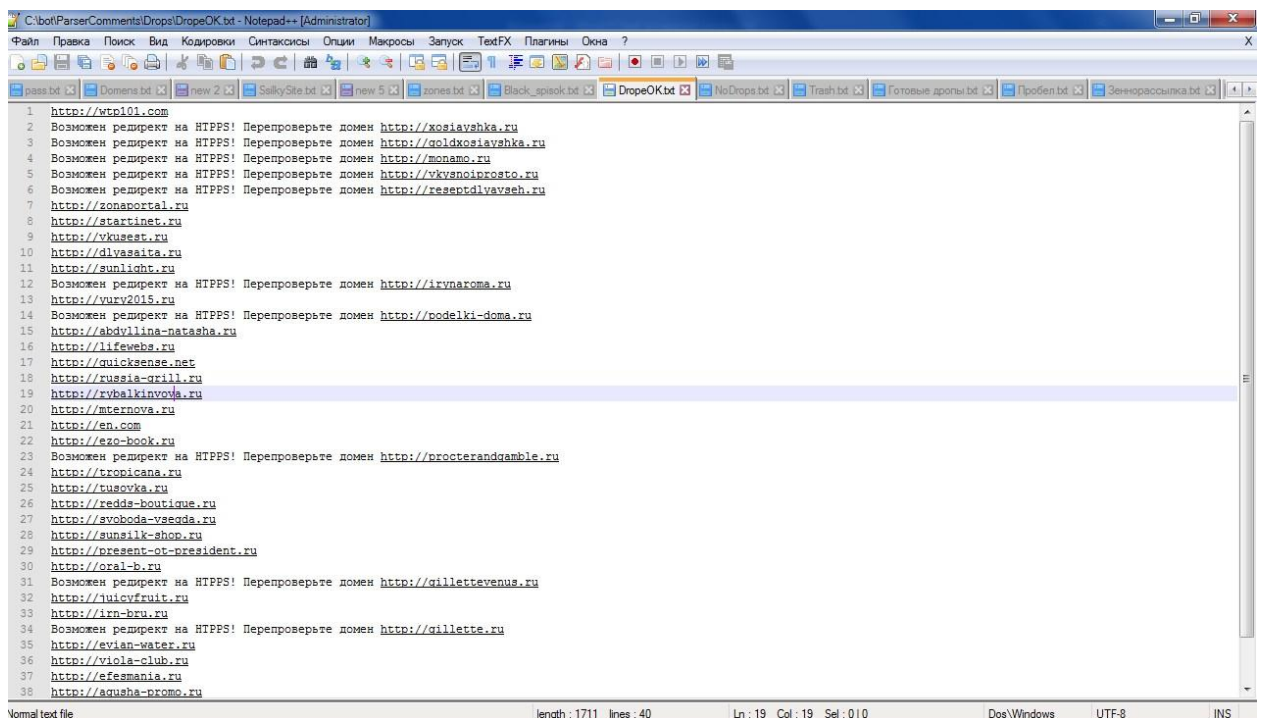
После этого скопировать УРЛы сайтов и записать их в файл **SsilkySite.txt**



```
1 http://neven.ru
2 http://ubavana.ru
3 http://astrologhiromant.ru
4 http://damskie-strasti.ru
5 http://zhenskajslabost.ru
6 http://schastliviyimir.ru
7 http://inseven.ru
8 http://qilaza.ru
9 http://promohunt.ru
10 http://silnet.ru
11 http://gallora.com
12 http://want2win.ru
13 http://lubimvidom.ru
14 http://academy.ru
15 http://kamagro.ru
16 http://avtoklav-store.ru
17 http://chiastoplotka.com
18 http://vashkofemem.ru
19 http://platform.ru
20 http://metabolismdiet.info
21 http://nistar.com
22 http://goiausa.info
23 http://beautypatrol.ru
24 http://knaasotavovsem.ru
```

1. В папке **Drops** находятся файлы **DropeOK.txt**, куда записываются найденные .

Пример файла DropeOK.txt



```
1 http://vtrp101.com
2 Возможен редирект на HTTPS! Перепроверьте домен http://xosiavshka.ru
3 Возможен редирект на HTTPS! Перепроверьте домен http://goldxosiavshka.ru
4 Возможен редирект на HTTPS! Перепроверьте домен http://monamo.ru
5 Возможен редирект на HTTPS! Перепроверьте домен http://vkvsnoiprostu.ru
6 Возможен редирект на HTTPS! Перепроверьте домен http://reaseptdivavseh.ru
7 http://zonaportal.ru
8 http://startinet.ru
9 http://vkusest.ru
10 http://diyasaite.ru
11 http://sunlight.ru
12 Возможен редирект на HTTPS! Перепроверьте домен http://rvmaroma.ru
13 http://vuvrv2015.ru
14 Возможен редирект на HTTPS! Перепроверьте домен http://podelki-doma.ru
15 http://abdvlina-natasha.ru
16 http://lifevehs.ru
17 http://quicksense.net
18 http://russia-grill.ru
19 http://rvbalkinvova.ru
20 http://mternova.ru
21 http://en.com
22 http://ezo-book.ru
23 Возможен редирект на HTTPS! Перепроверьте домен http://procterandgamble.ru
24 http://cropticana.ru
25 http://rusovka.ru
26 http://redde-boutique.ru
27 http://svoboda-vseoda.ru
28 http://sunlight-shop.ru
29 http://present-ot-president.ru
30 http://oral-b.ru
31 Возможен редирект на HTTPS! Перепроверьте домен http://gillettevenus.ru
32 http://quicyfruit.ru
33 http://irn-bru.ru
34 Возможен редирект на HTTPS! Перепроверьте домен http://gillette.ru
35 http://evian-water.ru
36 http://viola-club.ru
37 http://efesmania.ru
38 http://aquasha-promo.ru
```

За час работы найдено около 40 тематических дропа.

3. в ФАЙЛЕ **NoDrops.txt** находятся существующие сайты, в подавляющем большинстве вашей тематики, этот список используется для фильтрации, как и файл **Trash.txt**, но вы можете использовать их также для себя. Как пример, вы можете скопировать их в главный файл **SsillySite.txt**, чтобы не искать вручную, и программа пройдет по ним.

В **Black_spisok.txt** заносятся данные доменов, по которым бот уже проходил, чтобы не дублироваться. Бота нужно погонять не менее 2-3 часов, чтобы он набрал необходимые данные.