

# Поиск спамеров на основе анализа пользователей сервиса Youtube

Филатов В.

Научный руководитель:

Суворов В., EMC

Май 2013

# Введение

- Community mining
- Библиографические данные
- Социальные графы
- Алгоритмы классификации данных
- Выделение спамеров

# Описание задачи

Выделение спамеров и спам-контента на основе данных пользователей Youtube

- «спамер» - пользователь, который выкладывает порно и другой контент, запрещенный на Youtube, поисковый спамер
- «спам-фанат» - пользователь (или бот), который следит за контентом спамеров

Гипотеза: спамеры и спам-фанаты связаны в «социальную сеть» и они обладают похожими характеристиками.

# Постановка задачи

- Создание механизма выделения спамеров на основе данных о пользователях известного ресурса
- Визуализация результатов

# Элементы исследования

- Youtube API
- Взаимоотношения между пользователями: комментарии, подписчики, избранное, и т.п.
- Построение социальных графов пользователей сервиса
- Определение характеристик спамеров

# Этапы решения задачи

- Автоматизированный сбор данных
- Унификация данных
- Создание единой базы данных её наполнение
- Построение социального графа пользователей

Вершины – пользователи

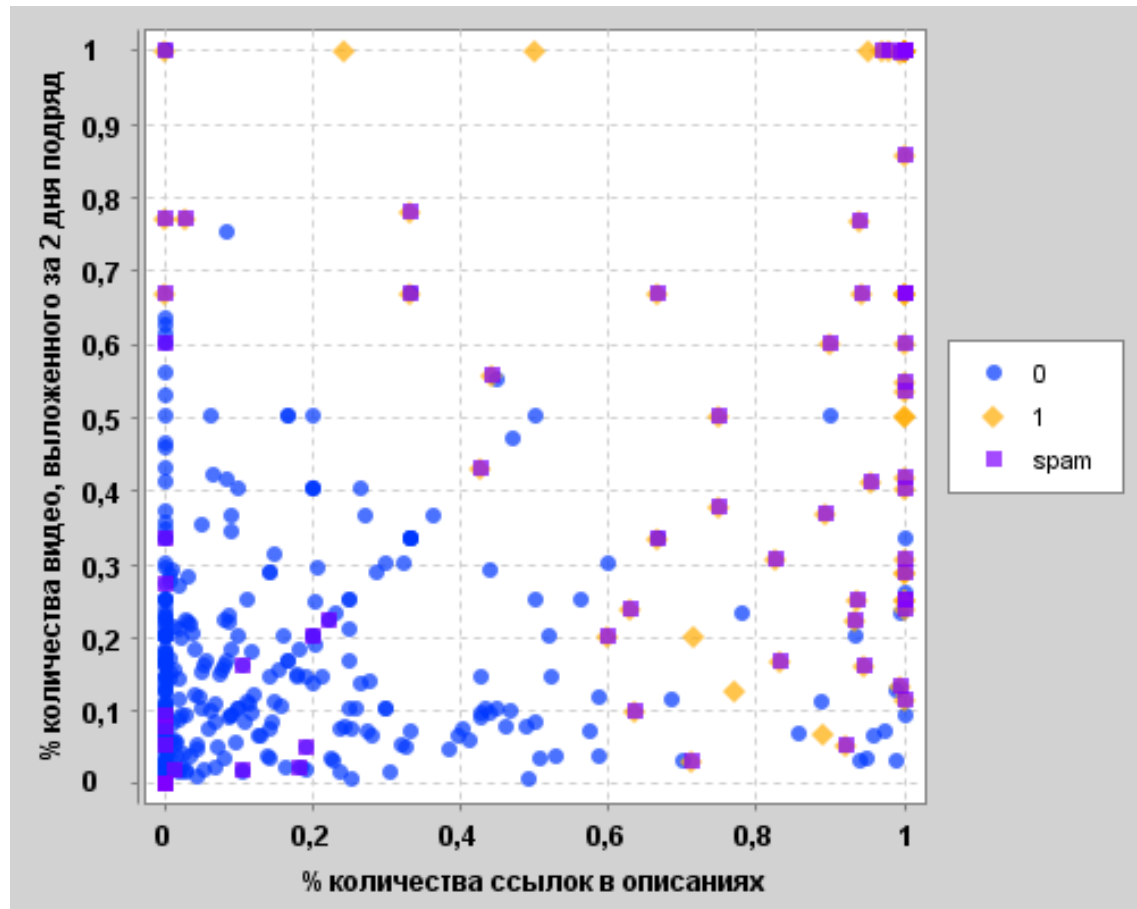
Связи – комментарии, подписчики,  
избранное

# Этапы решения задачи

- Выделение общих характеристик спамеров:
  - Ссылки в описаниях
  - Даты загрузки
  - Схожесть описаний
  - Активность пользователя
- Алгоритмы классификации данных
  - Kmeans
  - Байес
  - SVM

# Кластеризация

- Kmeans
- FP=0,15
- TN=0,15
- SVM
- Байес





# Социальная сеть

Firefox | Сообщение... | Введите Список Авторов | По

178.130.32.141:8084/finish/

Часто посещаемые

## Поиск спамеров на сервере YouTube

<< < Страница 1 из 23 > >> На странице 20 строк Установить

User	CountUploads	Cluters1	Clusters2	Count Subscribers	SameDescribes	degree	closeness	SameCommnt
<a href="#">theyoungturks</a>	14511	1	1	793425	0.0	19	1.0	1
<a href="#">moviemaniacsde</a>	4276	1	1	97103	0.16	1	1.0	2
<a href="#">collegehumor</a>	1934	1	1	4156550	0.0	29	1.0	5
<a href="#">psychetruth</a>	1697	1	1	238864	0.0	5	1.0	37
<a href="#">satanrulezzzz</a>	1436	1	1	33716	0.08	1	1.0	14
<a href="#">tytuniversity</a>	1379	1	1	82439	0.0	2	1.0	86
<a href="#">xtraonline</a>	717	1	1	6329	0.0	1	1.0	1
<a href="#">natgeowild</a>	709	1	1	152245	0.0	13	1.0	0
<a href="#">jpizzle1122</a>	666	0	1	604031	0.0	1	1.0	1351
<a href="#">hiphopnewzcentral</a>	530	1	1	1223	0.0	1	1.0	0
<a href="#">killuminaton</a>	486	1	2	92	0.0	1	0.0	2
<a href="#">gamebombcentral</a>	435	1	1	203391	0.64	8	1.0	28
<a href="#">mrchelos47</a>	324	1	1	594	0.0	1	1.0	0
<a href="#">thebreakingnews0</a>	315	1	1	1055	0.04	1	1.0	0
<a href="#">townsquare</a>	255	1	1	34325	0.0	2	1.0	2
<a href="#">microprikol</a>	201	1	1	165724	0.0	3	1.0	7
<a href="#">thisishoroshu</a>	200	1	1	1576224	0.0	3	1.0	0
<a href="#">ikrazeyyhd</a>	178	1	2	488	0.0	1	0.0	1
<a href="#">azatgw</a>	165	1	1	4884	0.24	1	1.0	54
<a href="#">miltonciyo</a>	158	1	1	1139	0.0	1	1.0	1

Автор skomeevskii  
Автор slavhygorkin  
Автор norozov1982  
Автор zajelimojnick

Добавить  
Send

# Результаты

Создан прототип веб-сервиса:

- Нахождение общих связей между пользователями сервиса
- Построение социального графа заданной группы пользователей

Выделение метрик пользователей сервиса

Алгоритм кластеризации пользователей

$FP=0.15$   $TN = 0.15$