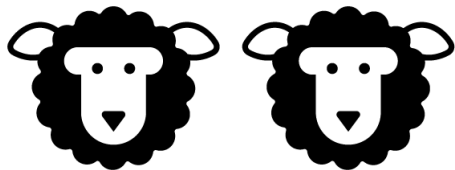


# Avito Duplicate Ads Detection



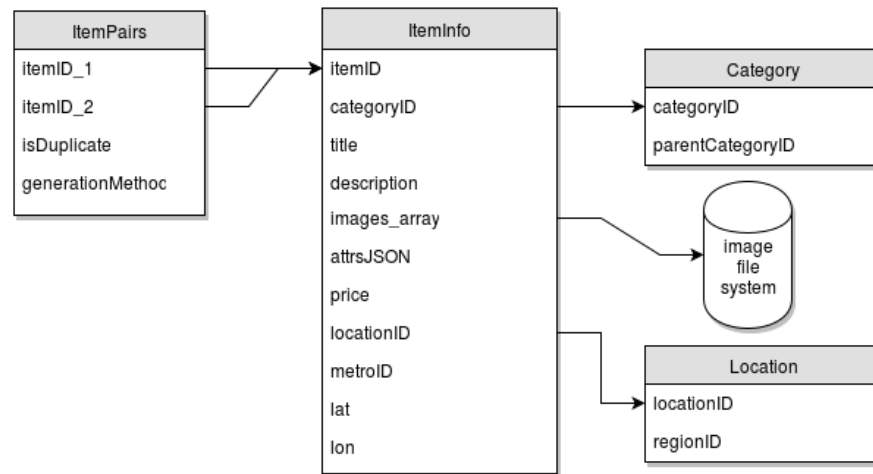
## AVITO DUPLICATE ADS DETECTION

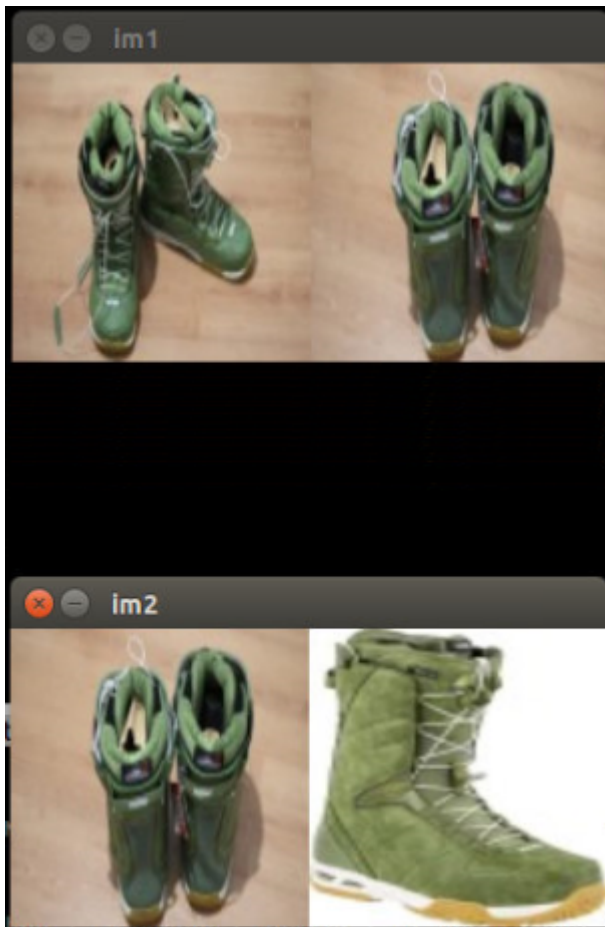
- Семенов Станислав , Цыбулевский Дмитрий, 2016

#	Δrank	Team Name <small>‡ model uploaded * in the money</small>	Score <small>?</small>	Entries	Last Submission UTC <small>(Best – Last Submission)</small>
1	—	<b>Devil Team <small>👤 *</small></b> <ul style="list-style-type: none"> <li>Stanislav Semenov</li> <li>u1234x1234</li> </ul>	<a href="#">0.95829</a>	162	<a href="#">Mon, 11 Jul 2016 23:09:03 (-0.7h)</a>
2	—	<b>TheQuants <small>👤 *</small></b> <ul style="list-style-type: none"> <li>Sonny Laskar</li> <li>NoName</li> <li>anokas</li> <li>Μαριος Μιχαηλιδης KazAnova</li> </ul>	<a href="#">0.95294</a>	197	<a href="#">Mon, 11 Jul 2016 19:53:15 (-46.7h)</a>
3	—	<b>ADAD <small>👤 *</small></b> <ul style="list-style-type: none"> <li>Gerard Toonstra</li> <li>Kele Xu</li> <li>Praveen Adepu</li> <li>Gilberto Titericz Junior</li> <li>Mario Filho</li> </ul>	<a href="#">0.94971</a>	226	<a href="#">Mon, 11 Jul 2016 23:57:54</a>
4	—	<b>8 + 9 = 11 <small>👤</small></b> <ul style="list-style-type: none"> <li>ZFTurbo</li> <li>Alexander Vikulin</li> </ul>	<a href="#">0.94694</a>	193	<a href="#">Mon, 11 Jul 2016 13:01:32 (-6.3h)</a>
5	—	<b>ololobhi <small>👤</small></b> <ul style="list-style-type: none"> <li>Abhishek</li> <li>ololo</li> </ul>	<a href="#">0.94587</a>	133	<a href="#">Mon, 11 Jul 2016 22:18:01</a>
6	—	<b>otivA</b>	<a href="#">0.94560</a>	117	<a href="#">Mon, 11 Jul 2016 13:09:48 (-0.1h)</a>
7	—	<b>Native Russian Speakers :P <small>👤</small></b> <ul style="list-style-type: none"> <li>Evgeny Eltyshev</li> <li>Georgiy Danshchin</li> </ul>	<a href="#">0.94449</a>	43	<a href="#">Mon, 11 Jul 2016 22:49:32 (-1.3h)</a>
8	<b>↑1</b>	<b>frist</b>	<a href="#">0.94438</a>	158	<a href="#">Mon, 11 Jul 2016 21:28:56 (-1h)</a>

# Данные

- 2991396 – тренировочная выборка
- 1044197 – тестовая выборка
- 10824317 Изображений





## Сноуборд ботинки Nitro Team 10 us

сноубордические ботинки Nitro Team

Размер 42,5

28см, 10 us

новые!!!

---

## Сноубордические ботинки

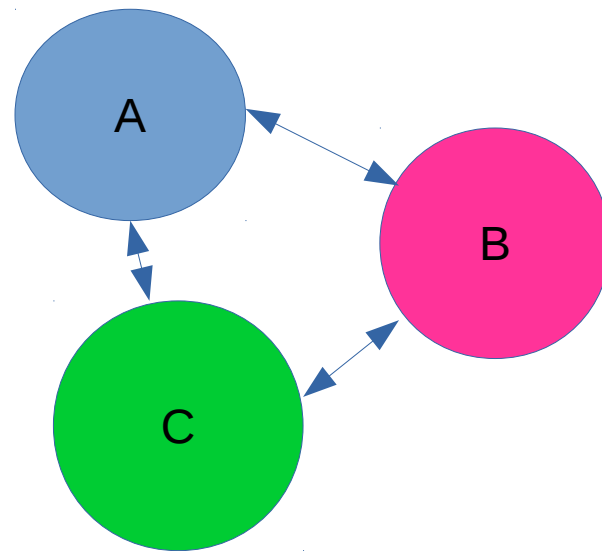
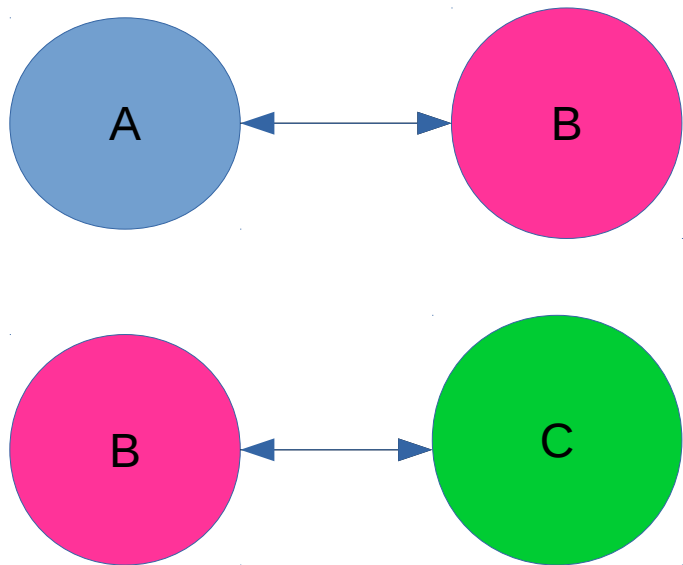
ботинки Nitro Team

Размер 42,5

28см, 10 us

новые!

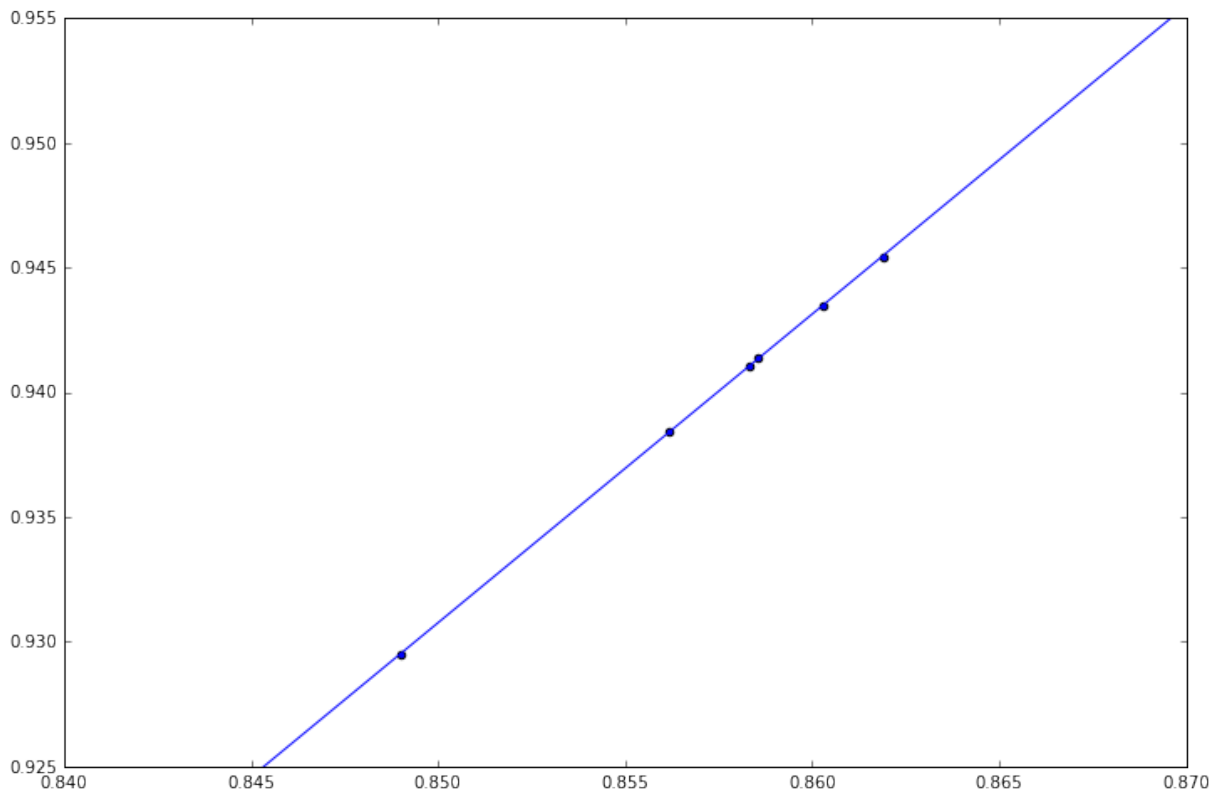
# Валидация



# Сайлент Режим

Заменяем каждое 10е значение на 0 (минимальное из всех).

По оси x – сайлент режим, по оси y – нормал.



Нужно заметить, что у нас объект – не просто объявление, а пара. Поэтому, не очень хорошо фичи по одному объявлению совать как просто `feature_ad1`, `feature_ad2`. Я всегда в одну колонку записывал минимум из этих двух, в другую – максимум из них. Т.е. например, получается минимум(цена двух объяв), максимум(цена двух объяв). В дальнейшем, если в качестве фичи указывается фича именно по одному объявлению, значит используется логика с мин+макс.

# Общие фишки

Цена, абсолютная разность цен, относительная разность цен, есть ли цены вообще, широта, долгота, декартово расстояние между городами, есть ли картинки вообще, совпадение широты, совпадение долготы, совпадение метро, совпадение региона, совпадение полностью всех json-аттрибутов, совпадение полностью тайтла, совпадение полностью дескрипшна и т.д.



# Числовые фичи

Вычленение всех чисел из тайтла, дескрипшена, тайтла+дескрипшена

Количество чисел в объявлении, количество повторяющихся чисел, относительное количество повторяющихся чисел, количество уникальных чисел, относительное количество уникальных чисел, сравнение медиан наборов чисел двух объявлений и т.д.

# Json фичи

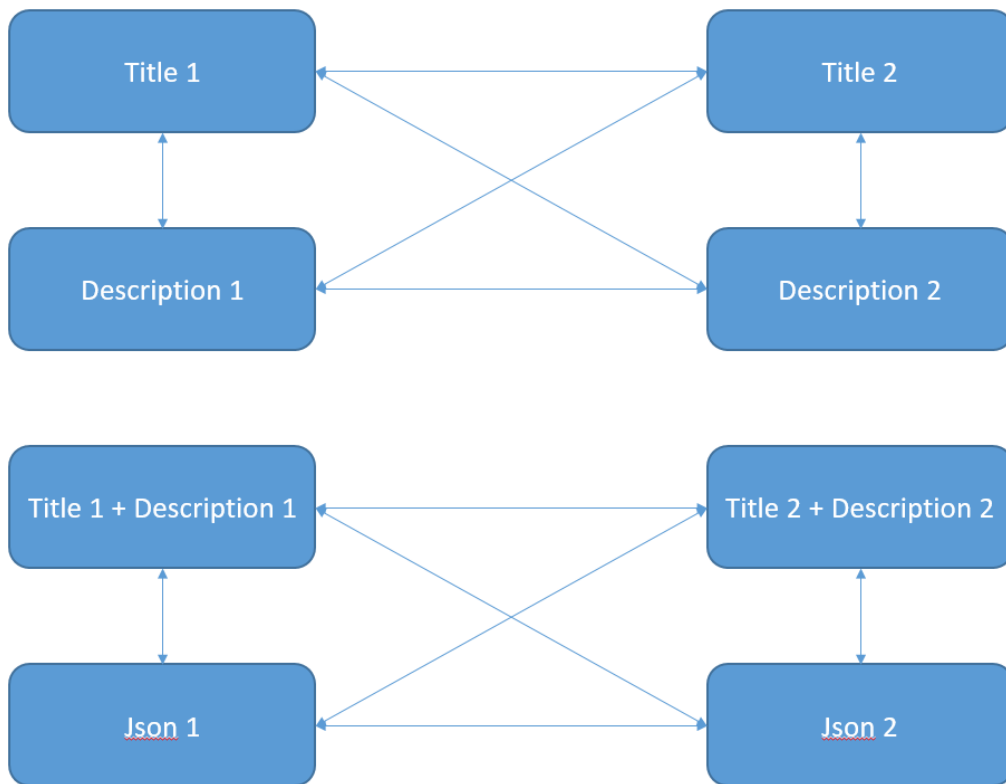
Общее количество json-аттрибутов  
объявлений, разность этого количества,  
относительная разность количества, число  
уникальных json-аттрибутов, число  
совпадающих json-аттрибутов, относительное  
число совпадающих json-аттрибутов и т.д.

# Текстовые фи́чи

Наборы: тайтл vs дескрипшн,  
тайтл + дескрипшн vs json

Сущность: слова, символы,  
триграммы символов

Косинусная мера между  
сущностями: тайтл-тайтл,  
дескрипшн-дескрипшн, тайтл-  
дескрипшн одного объявления,  
тайтл-дескрипшн разных  
объявлений, тайтл+дескрипшн-  
тайтл+дескрипшн, json-json,  
тайтл+дескрипшн-json одного  
объявления, тайтл+дескрипшн-  
json разных объявлений

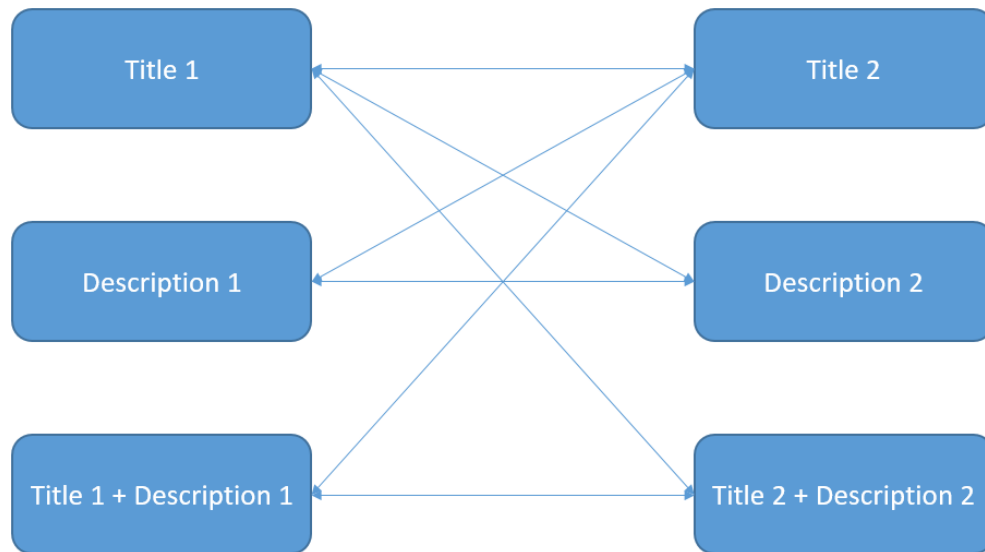


# Текстовые фи́чи

Наборы: тайтл, дескрипшн, тайтл + дескрипшн

Сущность: слова, символы

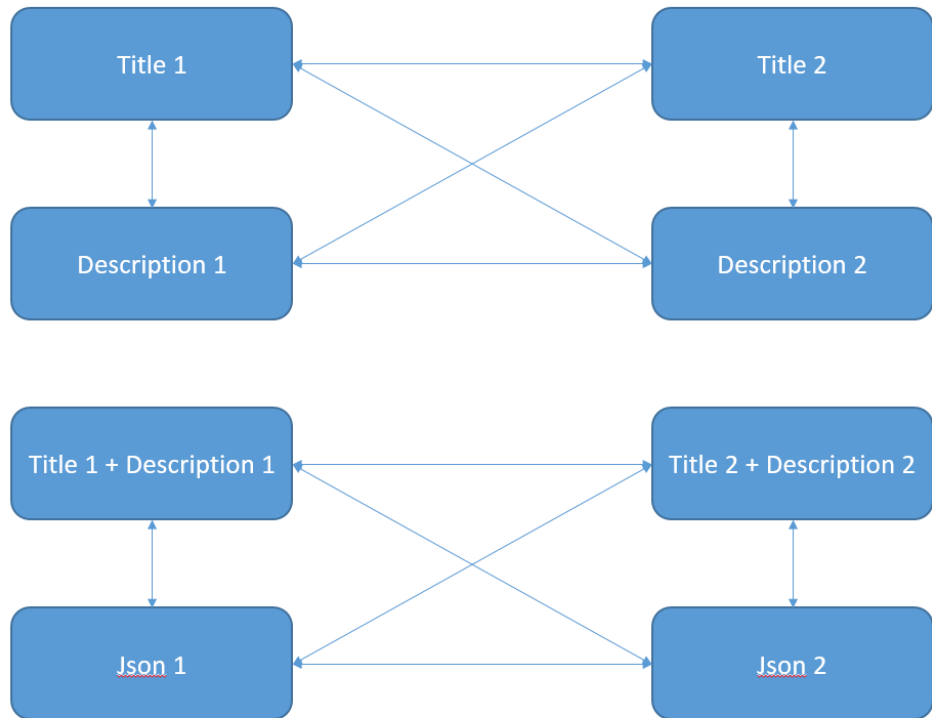
Число сущностей в каждом из наборов, относительное покрытие сущностей тайтла к сущностям тайтл+дескрипшн, абсолютная разность сущностей в каждом из наборов, относительная разность сущностей в каждом из наборов, матричное произведение сущностей каждого из наборов, матричное произведение между тайтлом и дескрипшном одного объявления, матричное произведение между тайтлом и дескрипшном разных объявлений, матричное произведение между тайтлом и тайтлом + дескрипшном разных объявлений, то же, но всё относительное и т.д.



# Word2Vec

Наборы: тайтл vs дескрипшн, тайтл + дескрипшн vs json

Обучаем ворд2век на тайтлах+дескрипшнах. Строим  $n\_similarity$  между: тайтл-тайтл, дескрипшн-дескрипшн, тайтл-дескрипшн одного объявления, тайтл-дескрипшн разных объявлений, тайтл+дескрипшн-тайтл+дескрипшн, json-json, тайтл+дескрипшн-json одного объявления, тайтл+дескрипшн-json разных объявлений



# LSI (LSA)

- На объединении объявлений
- На разности объявлений

# Текстовые фи́чи

- Различные расстояния title-title, title-description  
Jaccard, Levenshtein, Cosine, NCD
- Ручной парсинг единиц измерения (Gb)

# Текстовые фи́чи

- Stemming
- Lemmatization
- Transliteration
- Дешевый стемминг: `token = token[:-7]`



# Картиночные фи́чи

Совпадение хешей картинок при разных  
степенях сжатий

# Картиночные фиичи



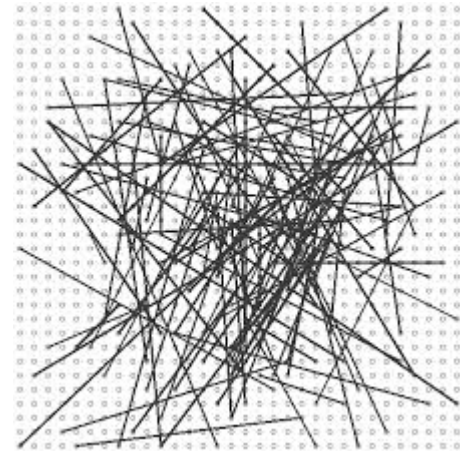
0.9, 0.5, 0, 0.8

Признаки - перцентили

# Картиночные фи́чи

- BRIEF: Binary Robust Independent Elementary Features

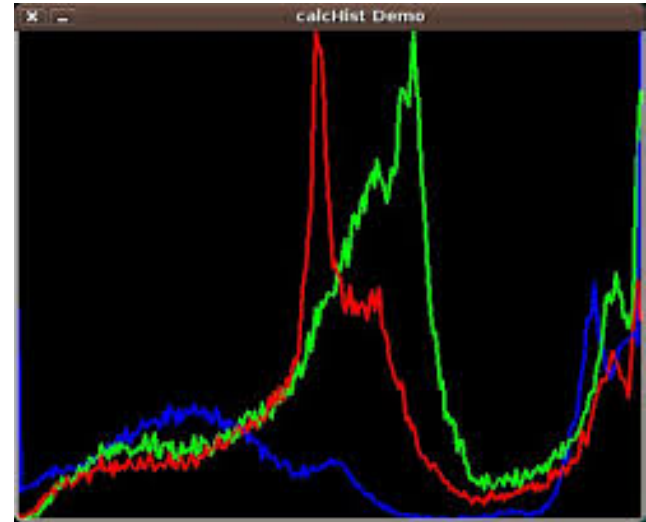
Hamming distance



# Картиночные фи́чи

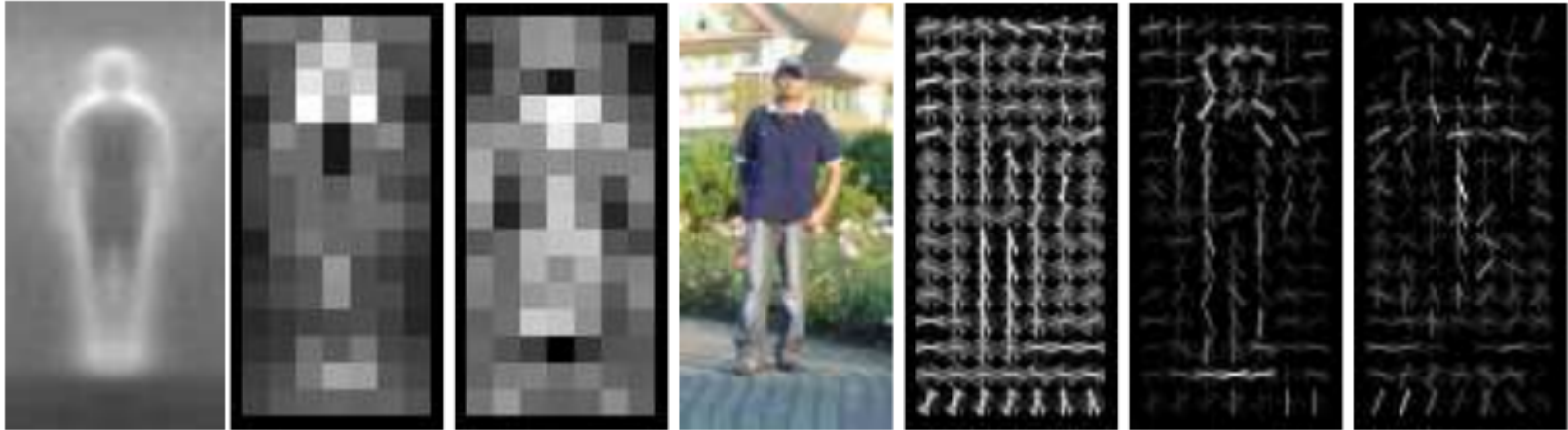
- Цветовые гистограммы

Пространство LAB, HSL  
Bhattacharyya distance



# Картиночные фи́чи

- Histogram of Oriented Gradients



# Картиночные фи́чи

- Нейронная Сеть обученная на ImageNet
- Full ImageNet Network
- Batch normalization: Accelerating deep network training by reducing internal covariate shift

# Картиночные фиичи

- AKAZE



# Построенные модели:

- 1) Различного вида NN
- 2) RF
- 3) ET
- 4) XGBoost
- 5) XGBoost без DevilFeas
- 6) XGBoost обученный только на generation\_method = 1
- 7) XGBoost обученный только на generation\_method = 3
- 8) XGBoost обученный только на generation\_method = 1 или 3
- 9) XGBoost обученный на транзитивной выборке (создание всех связей 1го уровня)
- 10) Xgboost – мультикласс, предсказывающий тип generation\_method (3 класса)
- 11) Xgboost – мультикласс, предсказывающий generation\_method-isDuplicate (5 классов)

Что дало основной вклад?

XGBoost, XGBoost обученный только на generation\_method = 1 или 3

Модель поверх них - XGBoost



# Обучение

- 7 серверов 40 / 160