



Introduction to **Active Learning**

Alexey Voropaev

go.mail.ru

A white hand cursor icon with a black outline, pointing upwards and slightly to the right, positioned over the "go.mail.ru" text.

Mail.Ru have own Search Engine



- It is not «patched Google»
- Completely out product
- 8.3% of market
- The same engine for:
 - web
 - image
 - video
 - news
 - real time
 - e-mail search
 - etc.

Mail.Ru Почта 2 Мой Мир Одноклассники Игры Знакомства Новости Поиск Все проекты ▾

поиск@mail.ru хью лори фото Найти


Интернет

- Картинки
- Видео
- Новости
- Обсуждения
- Ответы
- Словари

Весь интернет

Поиск в Москве

Картинки




Хью Лори (Hugh Laurie) фото

Если у вас есть интересные фото Хью Лори(Hugh Laurie) . Качественные фото Хью Лори - Hugh Laurie новости - Хью Лори хочет ездить по Москве на мотоцикле

theplace.ru/photos/Hugh_Laurie-... копия еще с сайта

Лори, Хью — Википедия



Джеймс Хью Кэлам Лори ОВЕ (англ. James Hugh Calum Laurie; род. 11 июня 1959 года в Оксфорде, Великобритания) — английский актёр, профессиональный пианист, сценарист, режиссёр, продюсер, писатель и певец, известность которому принесли прежде всего роли в британских комедийных телесериалах «Чёрная гадюка...»

ru.wikipedia.org/wiki/Лори,_Хью копия еще с сайта

Хью Лори : Фотографии

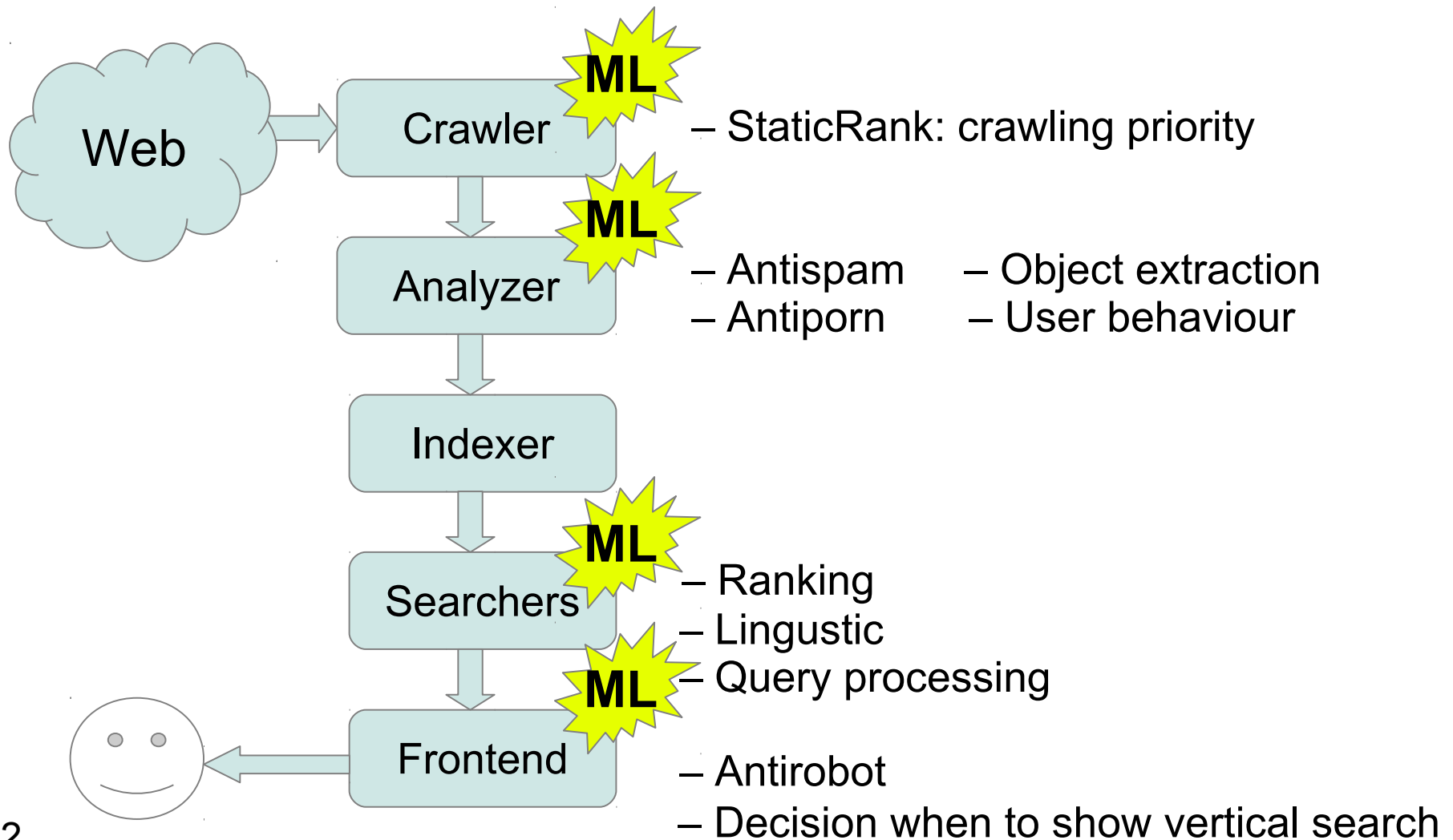
Хью Лори Hugh Laurie - Биография - Фильмография - Фото - Обои - Кадры - Награды - Новости - Отзывы - 500 x 667, 53kb - 500 x 581, 62kb - 500 x 626, 38kb

kinomania.ru/stars/h/Hugh_Lauri... копия еще с сайта

Хью Лори (Hugh Laurie) - биография и фото Хью Лори (Hugh Laurie...

Подпись имя: Хью Лори - Дата рождения: 11.06.1959 -

Simplified Search Engine Architecture:



«Machine learning is the hot new thing»

John Hennessy, President, Stanford



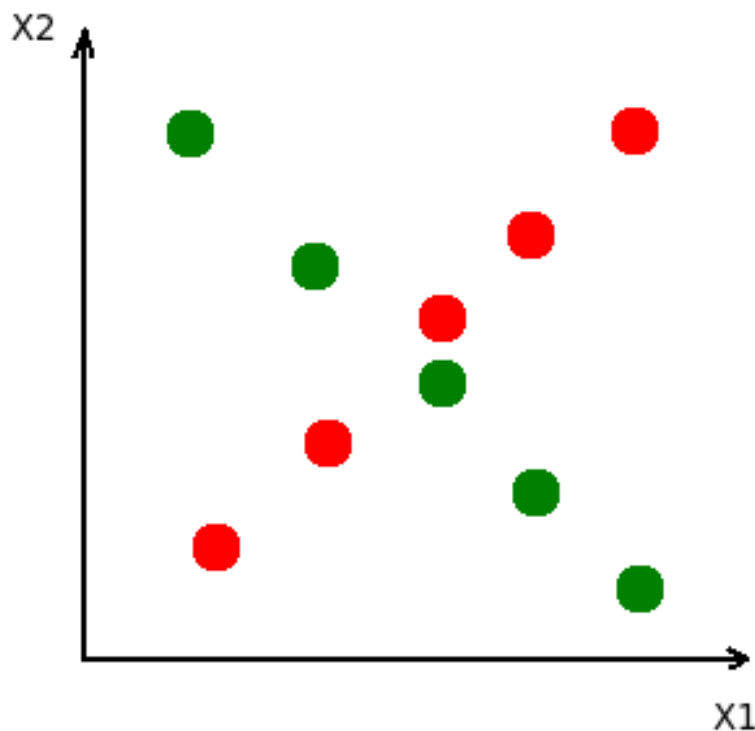
There is unknown function $f(\vec{x}) = y$

Given:

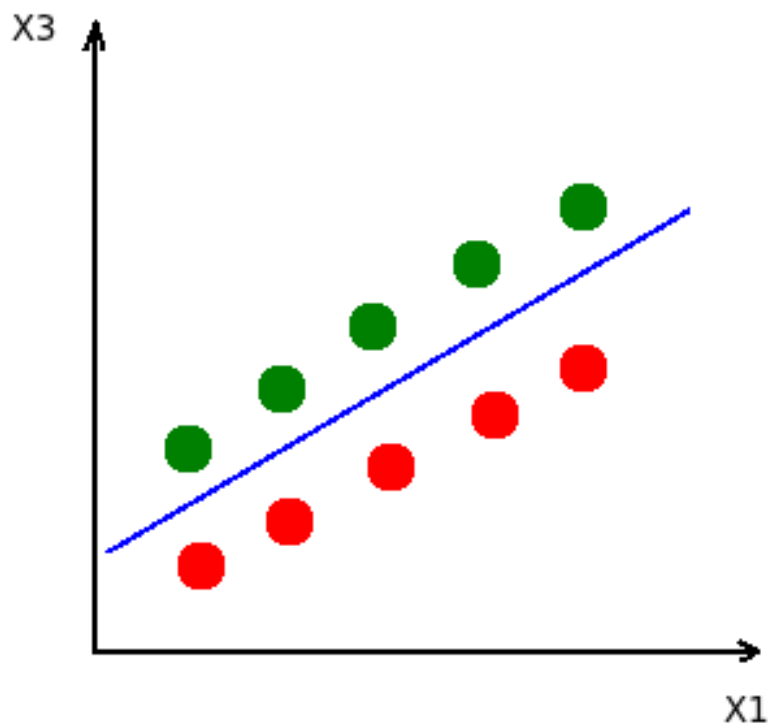
Set of samples $T = \{\vec{x}_1, y_1\} \dots \{\vec{x}_n, y_n\}$

Goal:

Build approximation f' of f using T



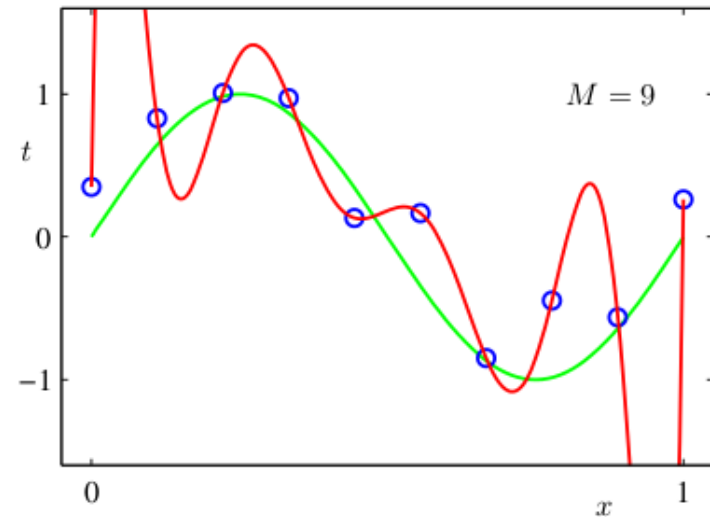
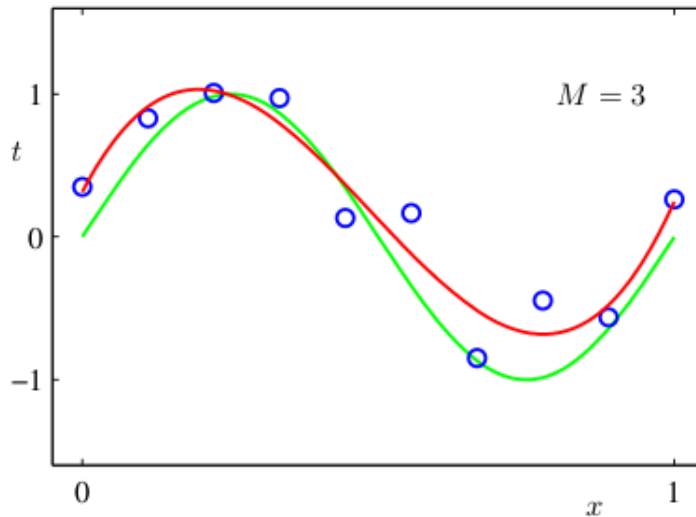
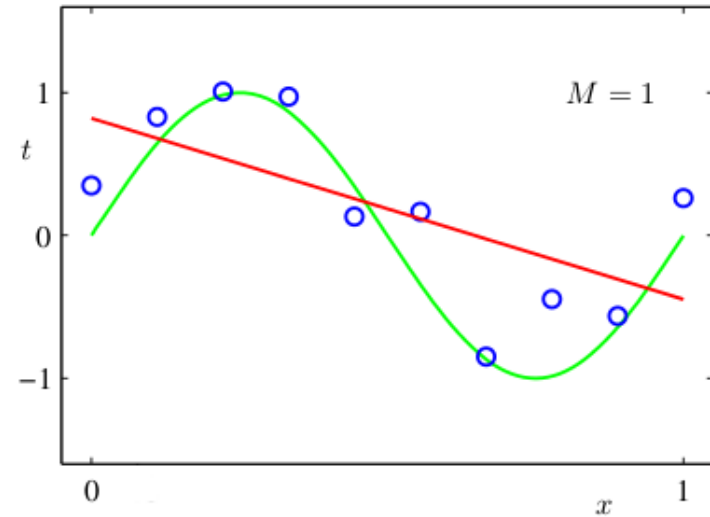
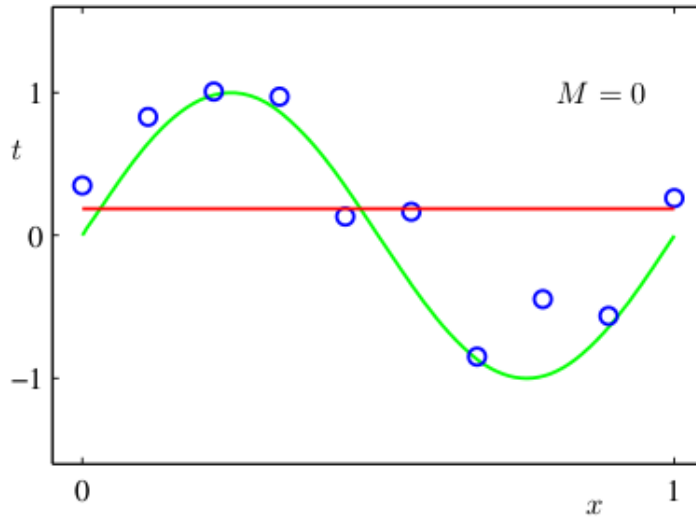
A — impossible

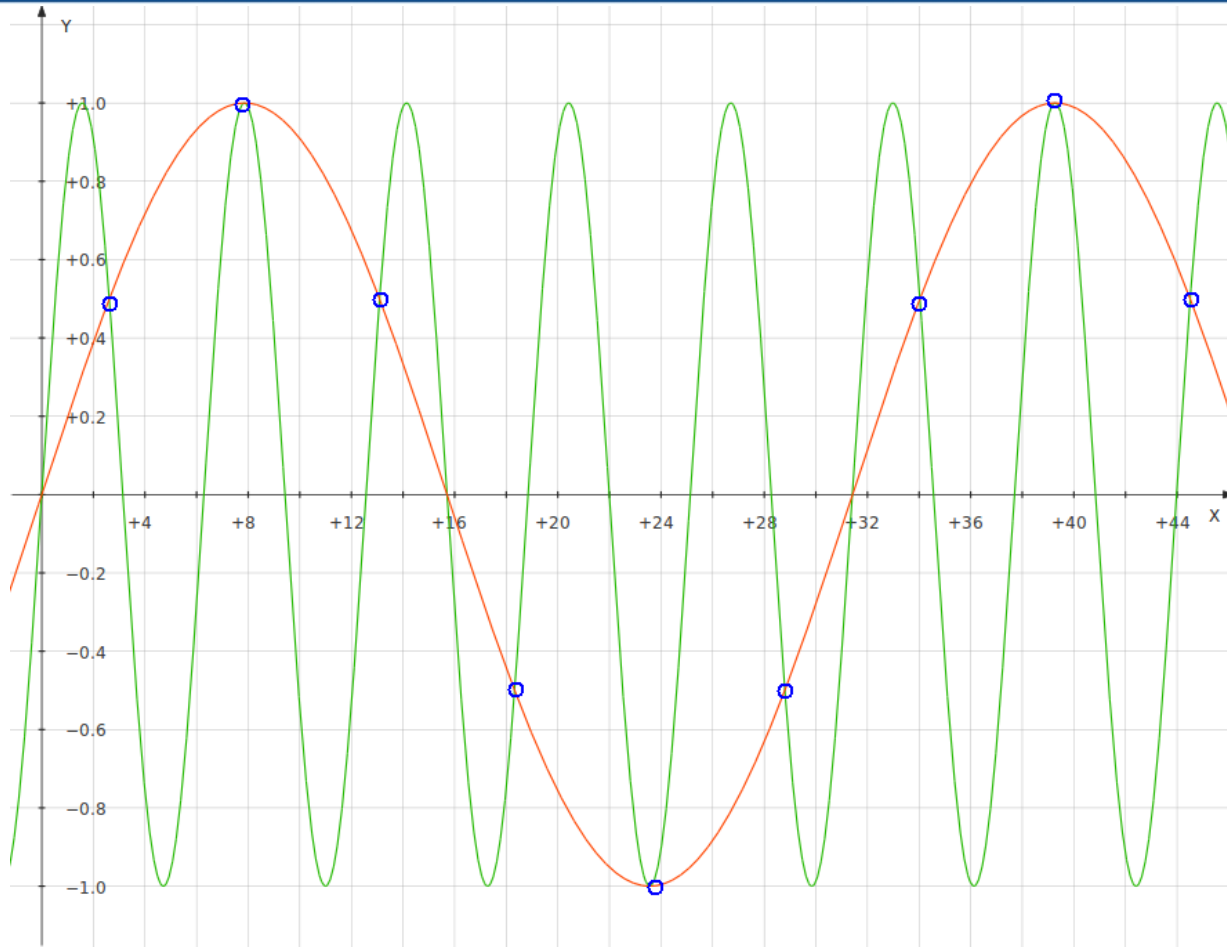


B — possible

Linear model

Proper model selection





Usually:

- random sampling is not the best strategy
- labeling is very expensive

Modern books about ML:

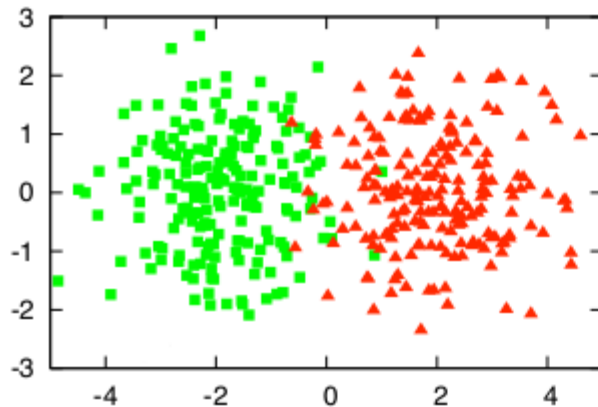
- Pattern Recognition and Machine Learning — 0
- Elements of Statistical Learning — 0
- An Introduction to Information Retrieval — 1 paragraph

Book about TS construction:

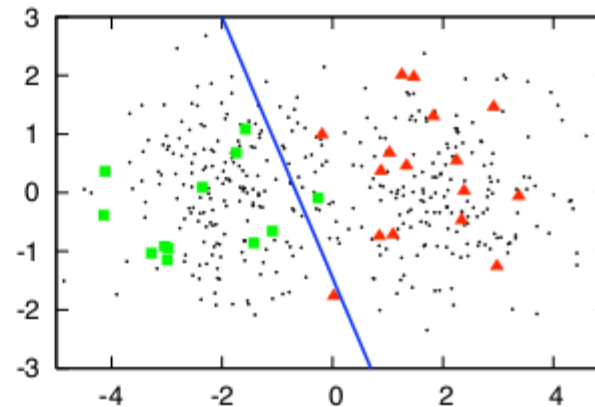
Burr Settles. *Active learning*.

- Publication Date: 2 July 2012

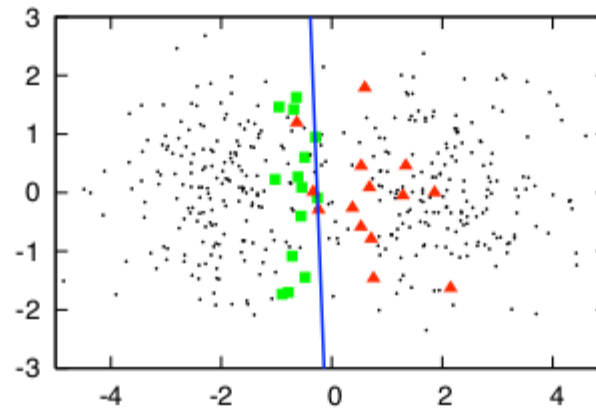
Idea of Active Learning



(a)

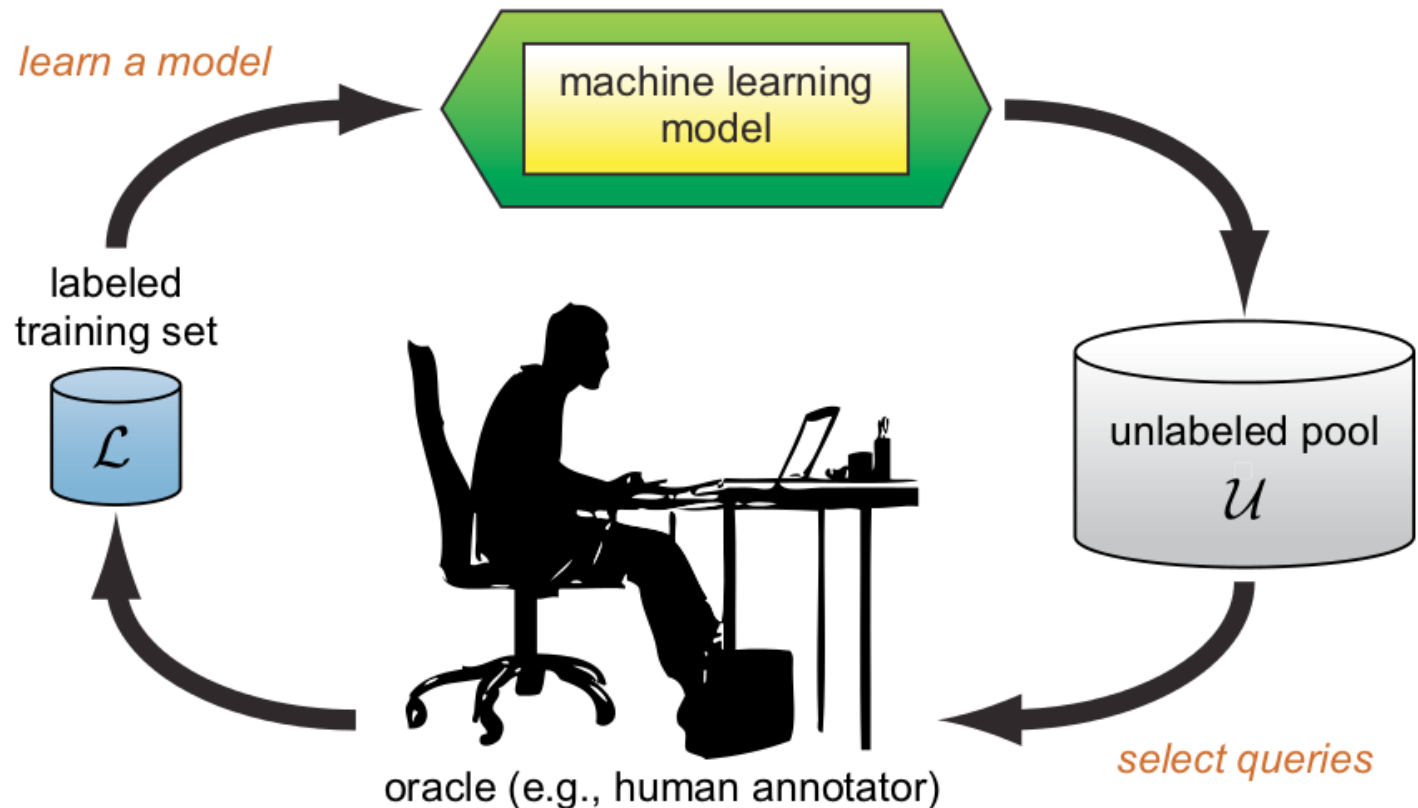


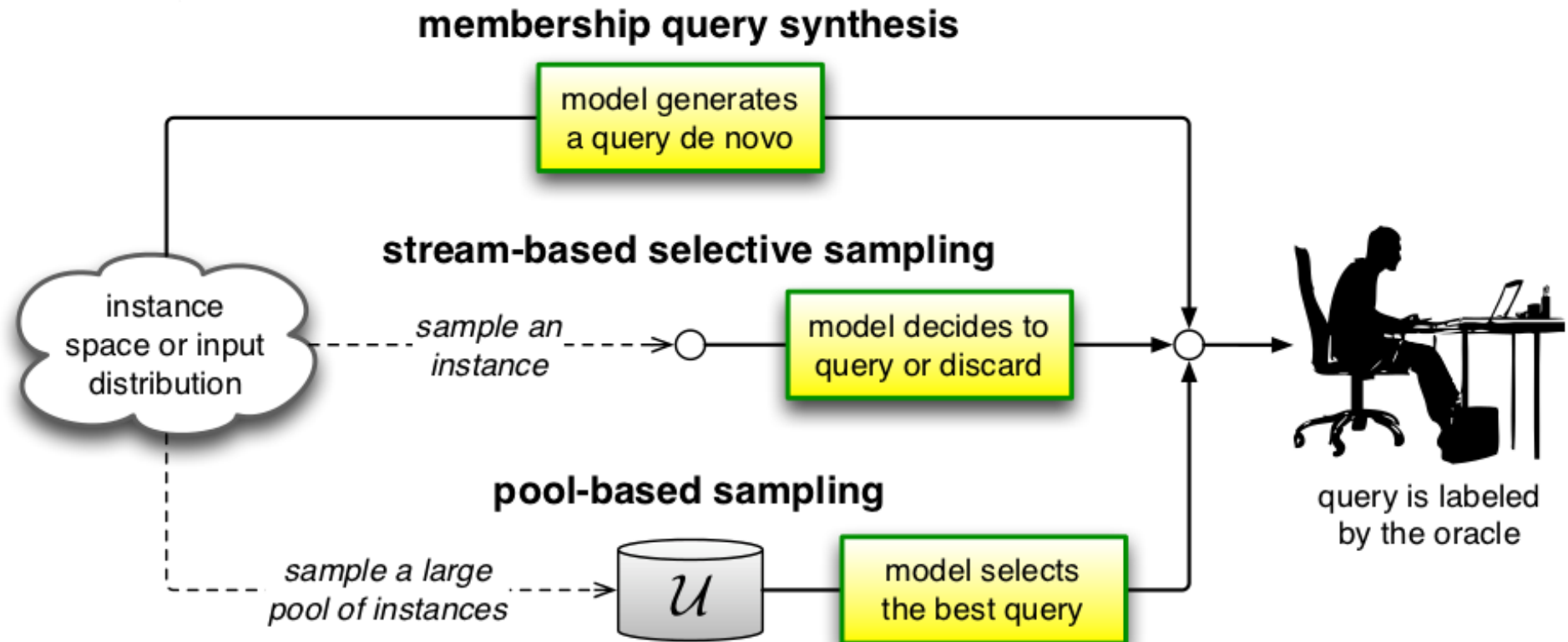
(b)



(c)

Main assumption: obtaining an unlabeled instance is free





Take instances about which it is least certain how to label.

Least confident:

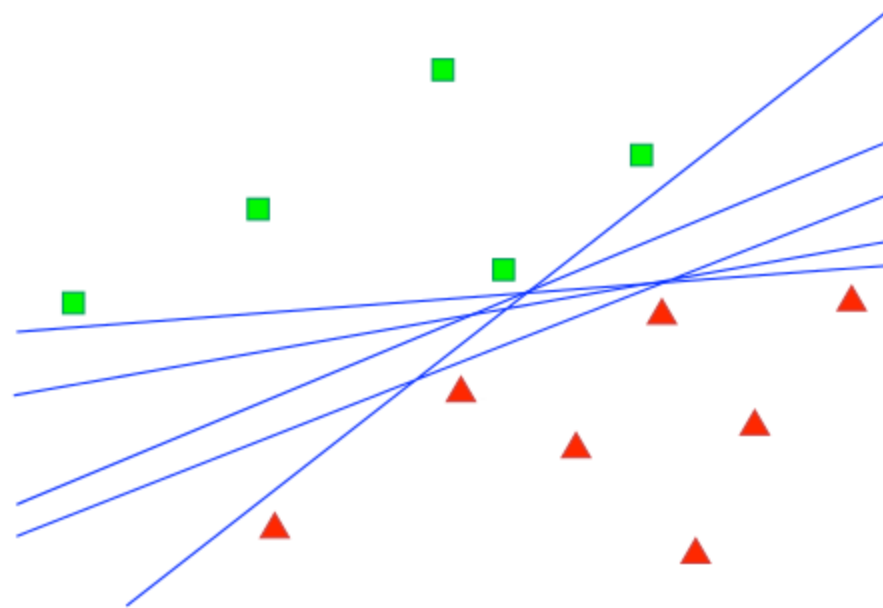
$$x_{LC}^* = \operatorname{argmax}_x 1 - P_{\theta}(\hat{y}|x)$$

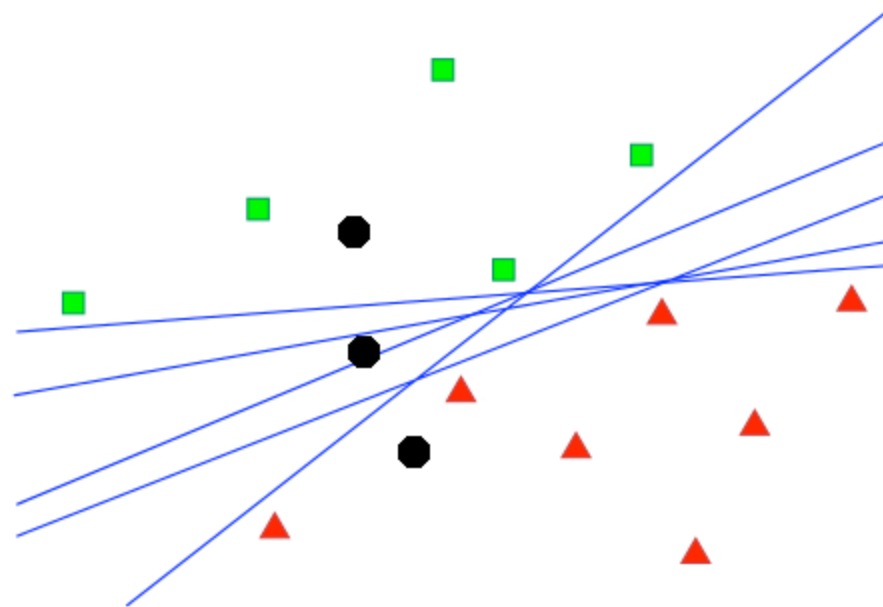
Margin sampling

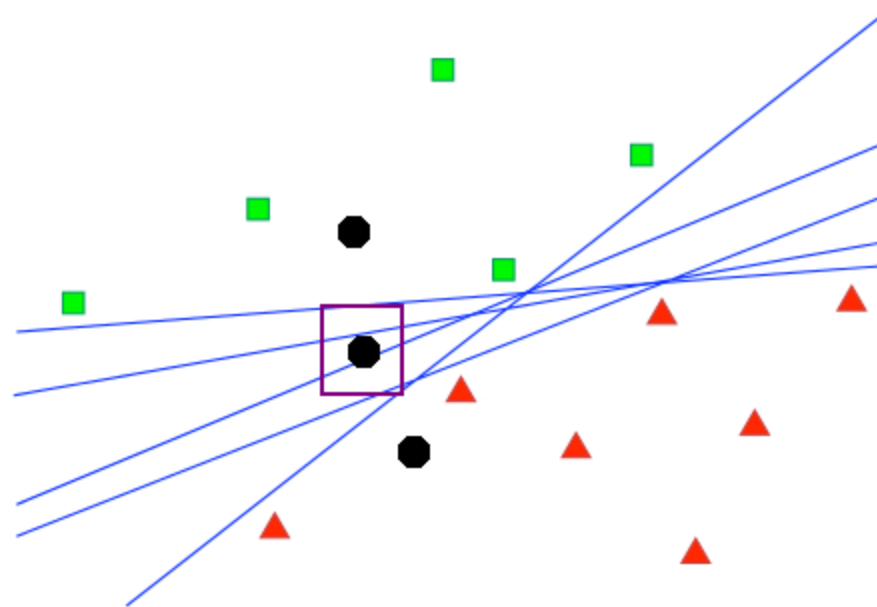
$$x_M^* = \operatorname{argmin}_x P_{\theta}(\hat{y}_1|x) - P_{\theta}(\hat{y}_2|x)$$

Entropy

$$x_H^* = \operatorname{argmax}_x - \sum_i P_{\theta}(y_i|x) \log P_{\theta}(y_i|x)$$







Measure of committee disagreement:

Vote entropy

$$x_{VE}^* = \operatorname{argmax}_x - \sum_i \frac{V(y_i)}{C} \log \frac{V(y_i)}{C}$$

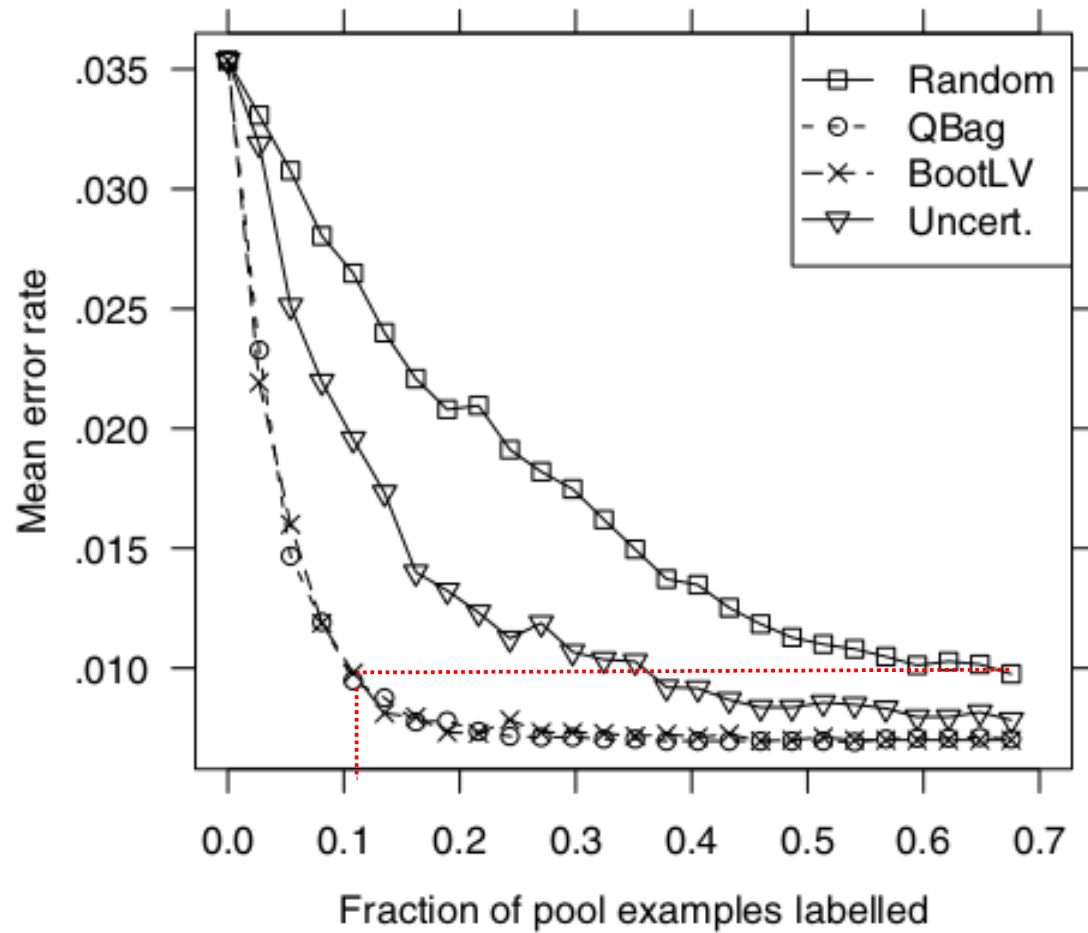
*Kullback-Leibler (KL)
divergence:*

$$x_{KL}^* = \operatorname{argmax}_x \frac{1}{C} \sum_{c=1}^C D(P_{\theta(c)} \| P_C),$$

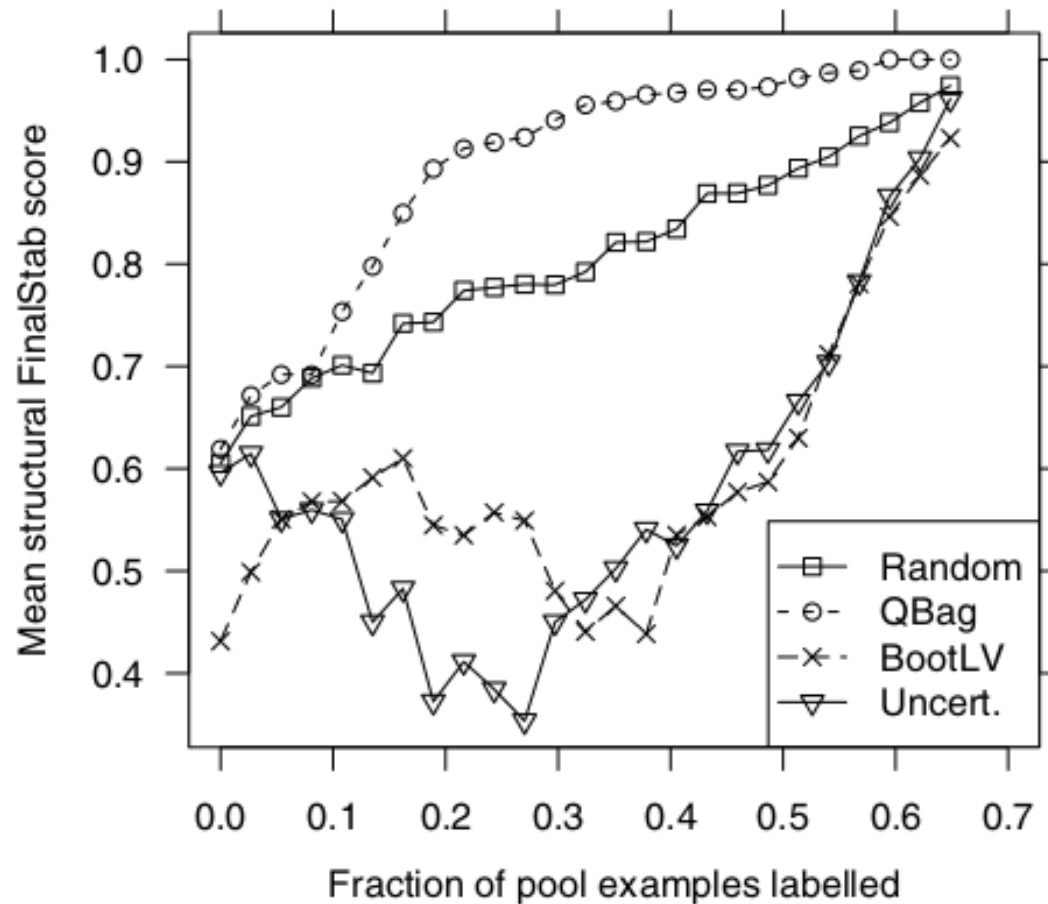
$$D(P_{\theta(c)} \| P_C) = \sum_i P_{\theta(c)}(y_i|x) \log \frac{P_{\theta(c)}(y_i|x)}{P_C(y_i|x)}$$

Input: T – labeled train set
 C – size of the committee
 A – learning algorithm
 U – set of unlabeled objects

1. Uniformly resample T , obtain $T_1 \dots T_C$, where $|T_i| < |T|$
2. For each T_i build model H_i using A
3. Select $x^* = \min_{x \in U} | |H_i(x)=1| - |H_i(x)=0| |$
4. Pass x^* to oracle and update T
5. Repeat from 1 until convergence

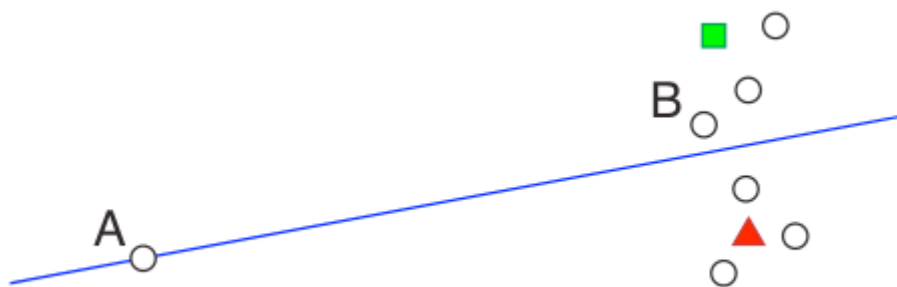


(b) Error rate



(a) Structural FinalStab ($\epsilon = 0$)

Idea: Inhabit dense/sparse regions of the input space



$$x_{ID}^* = \operatorname{argmax}_x \phi_A(x) \times \left(\frac{1}{U} \sum_{u=1}^U \operatorname{sim}(x, x^{(u)}) \right)^\beta$$

Expected Model Change

Expected Error Reduction

Variance Reduction

Many «model dependent» methods

Example 1: Recognition of sentence boundaries



«I recommend Mitchell (1997) or Duda et al. (2001). I have strived...»



Idea:

1. Make set of simple regexp-based rules like
«*big letter at the right*» or
«*digits at the right and left*»
(40 rules)
2. Build classifier:
Gradient Boosted Decision Trees

A.Kudinov, A.Voropaev, A.Kalinin. A HIGHLY PRECISIONAL METHOD FOR THE RECOGNITION OF SENTENCE BOUNDARIES. Computational Linguistics and Intellectual Technologies. 2011

Example 1: Recognition of sentence boundaries



OpenNLP /
Sentence
Boundary
Detector /
Random
sampling

«AOT»
project

Mail.Ru
detector /
Random
sampling

Mail.Ru
detector /
least
confident

Error
Rate

41,2 %

30,4 %

8.2 %

0,8 %

Train set: 9820 examples

Validation: 500 examples

Example 2: Ranking formula

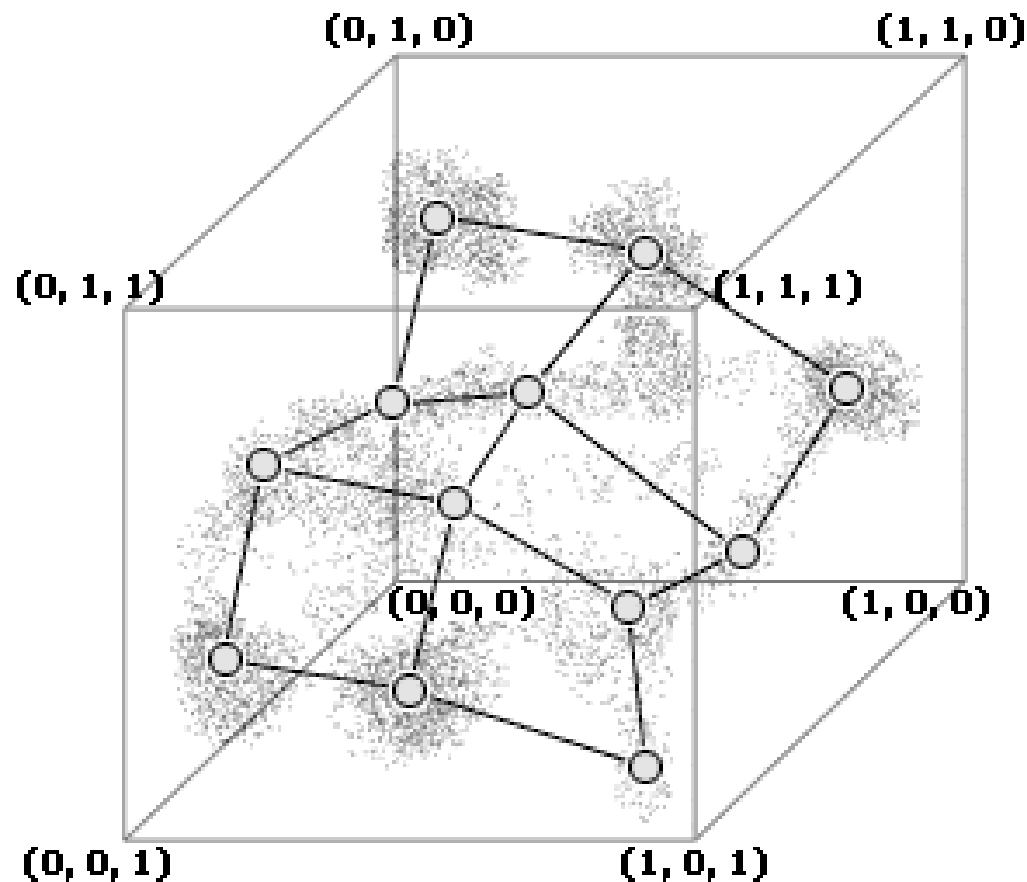


Idea:

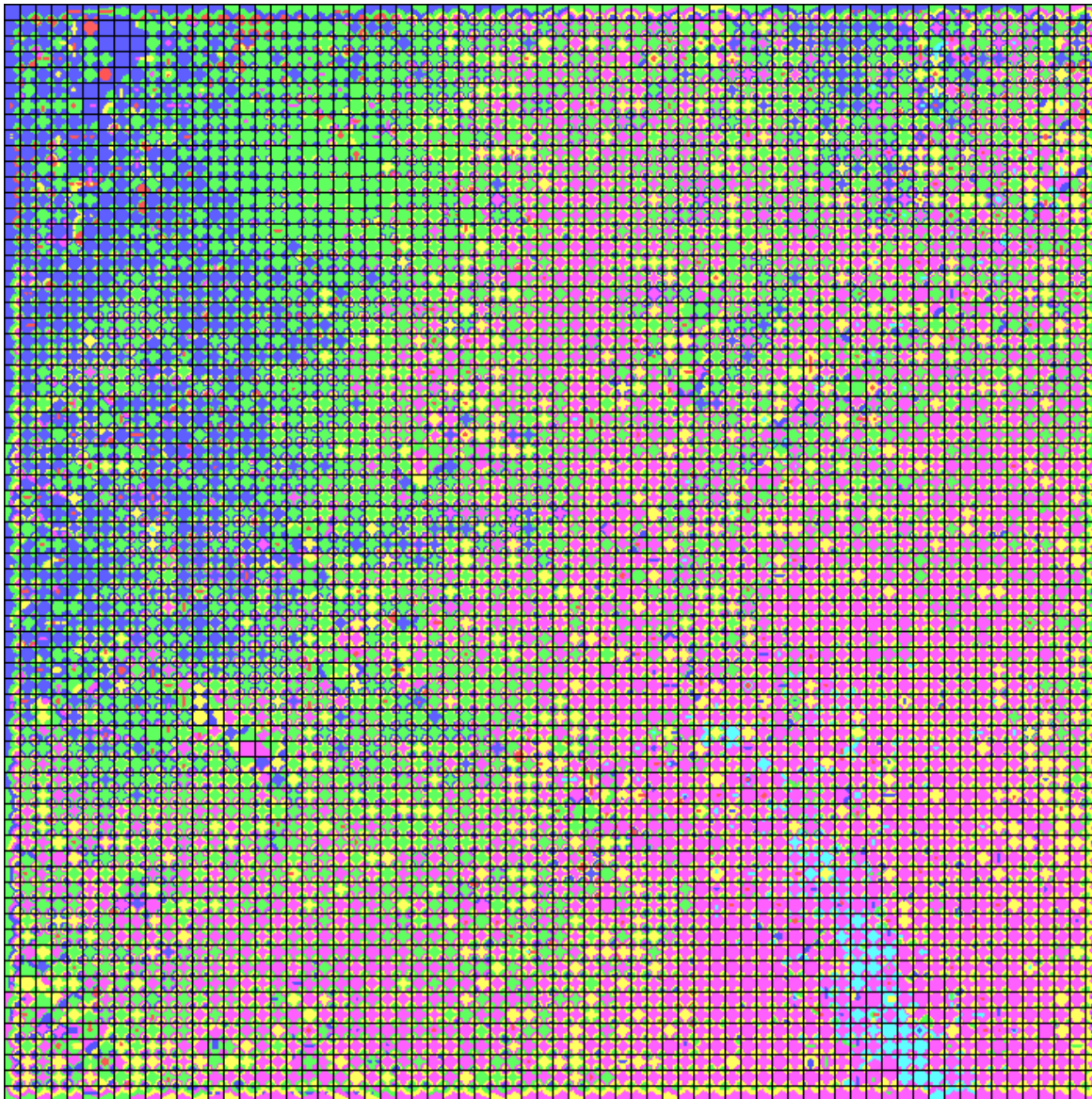
1. Query-Document presented by x , where:
 - x_1 – tf-idf rank
 - x_2 – query geographical
 - x_3 – geo-region of user is equal to document's region
(600 other factors)
2. Build train set:
 - 5 – vital
 - 4 – exact answer
 - 3 – usefull
 - 2 – slightly usefull
 - 1 – out of topic
 - 0 – can't be labeled (unknown language, etc.)
3. Train ranking formula using modified LambdaRank

Example 2: Self-organizing map

Idea: mapping from N-dimensional to 2-dimensional

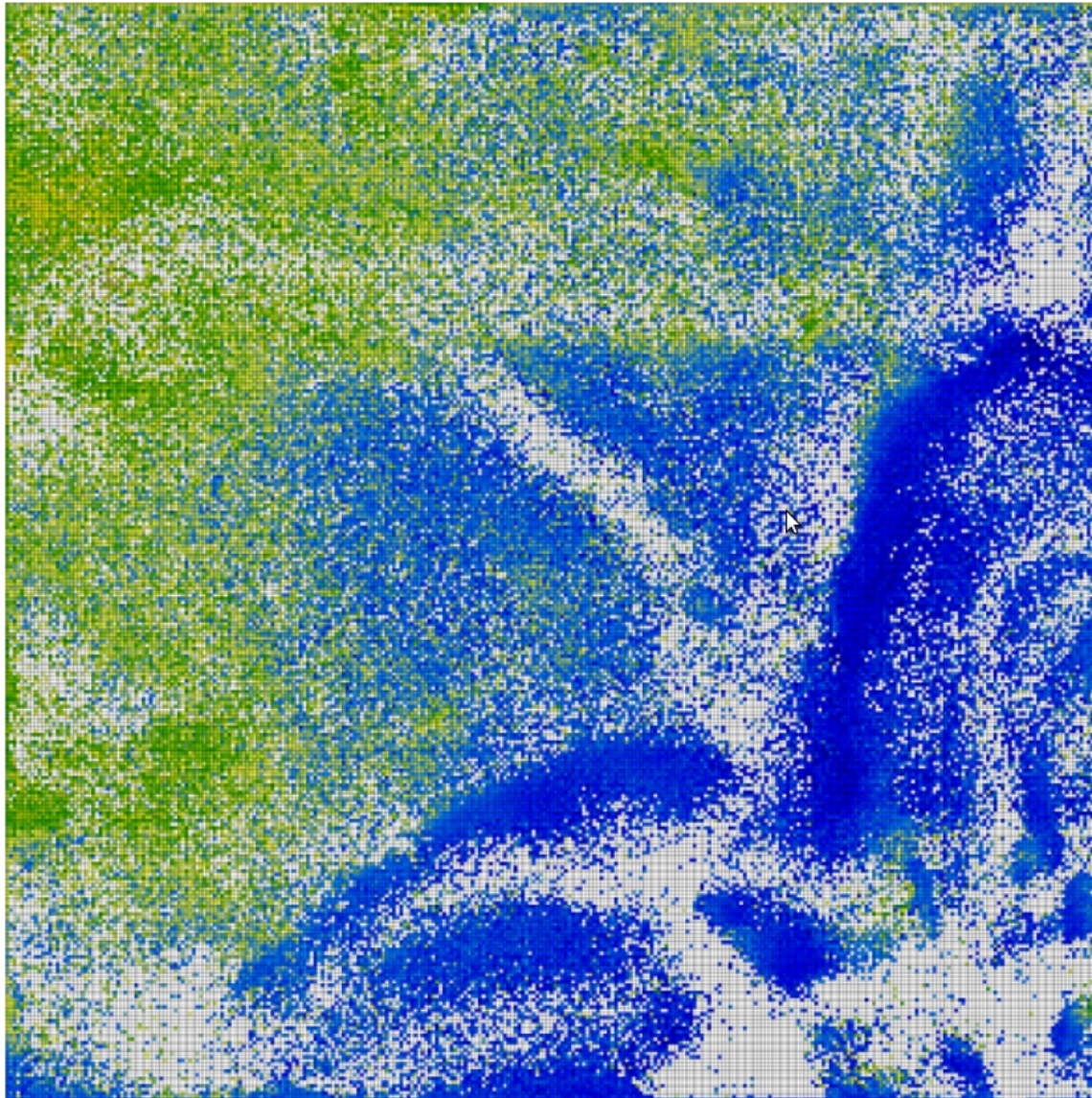


Example 2: Map of train set for ranking



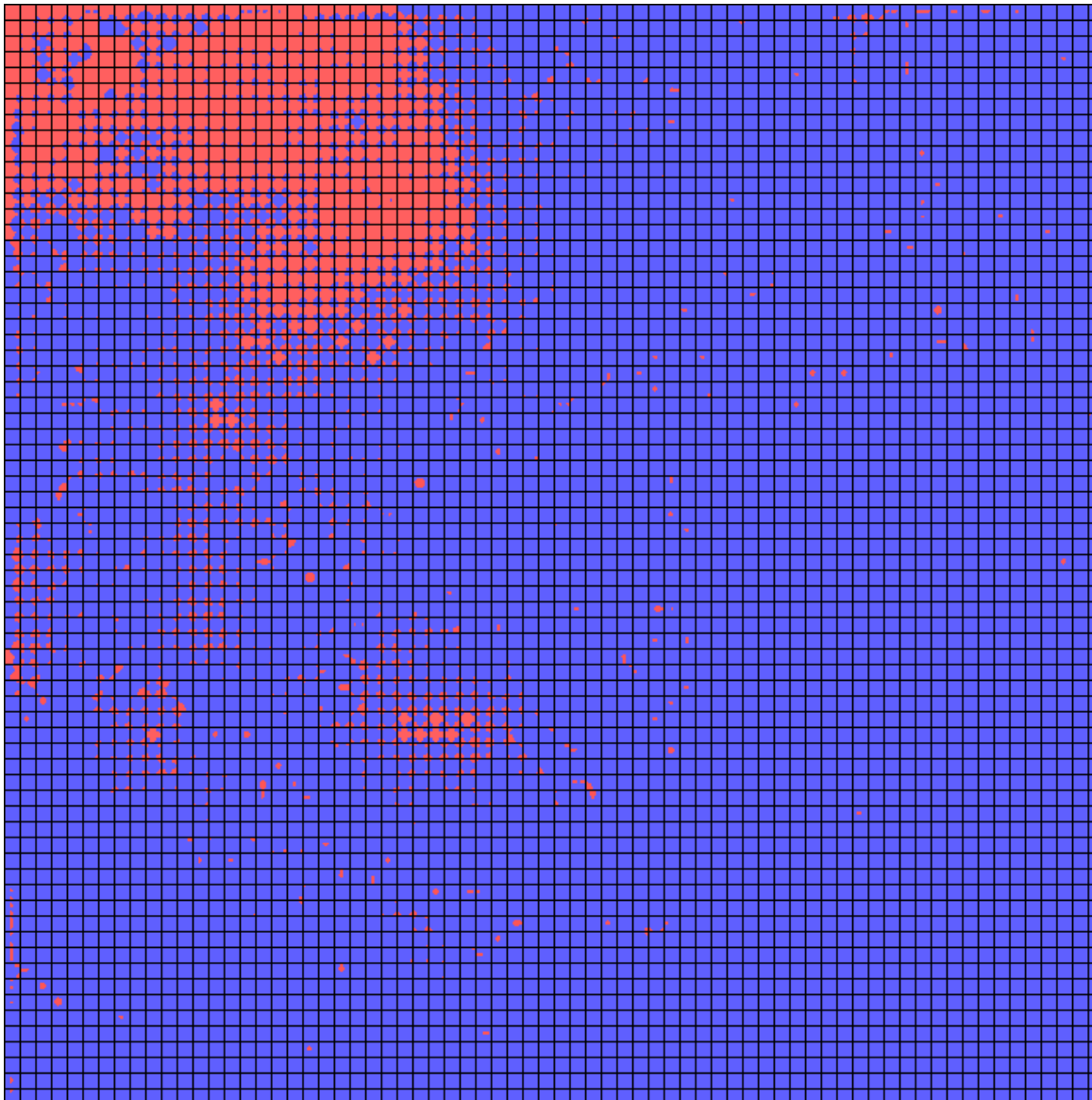
eval_0	
eval_1	
eval_2	
eval_3	
eval_4	
eval_5	

Example 2: Density of QDocuments on the map



High density
Low density

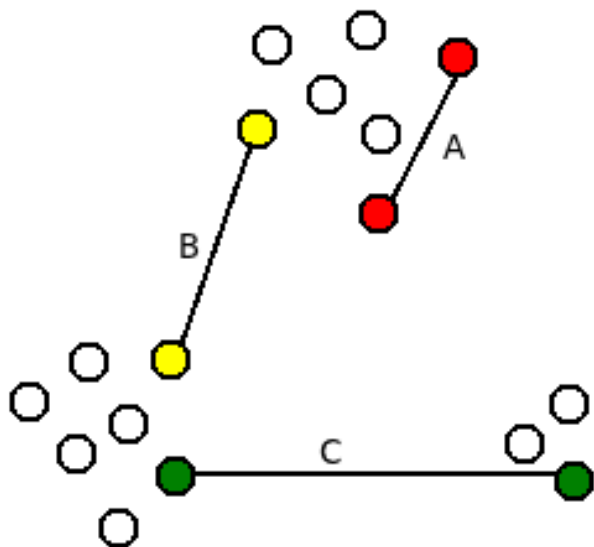
Example 2: Result of “SOM miss” selection



Old documents
New documents

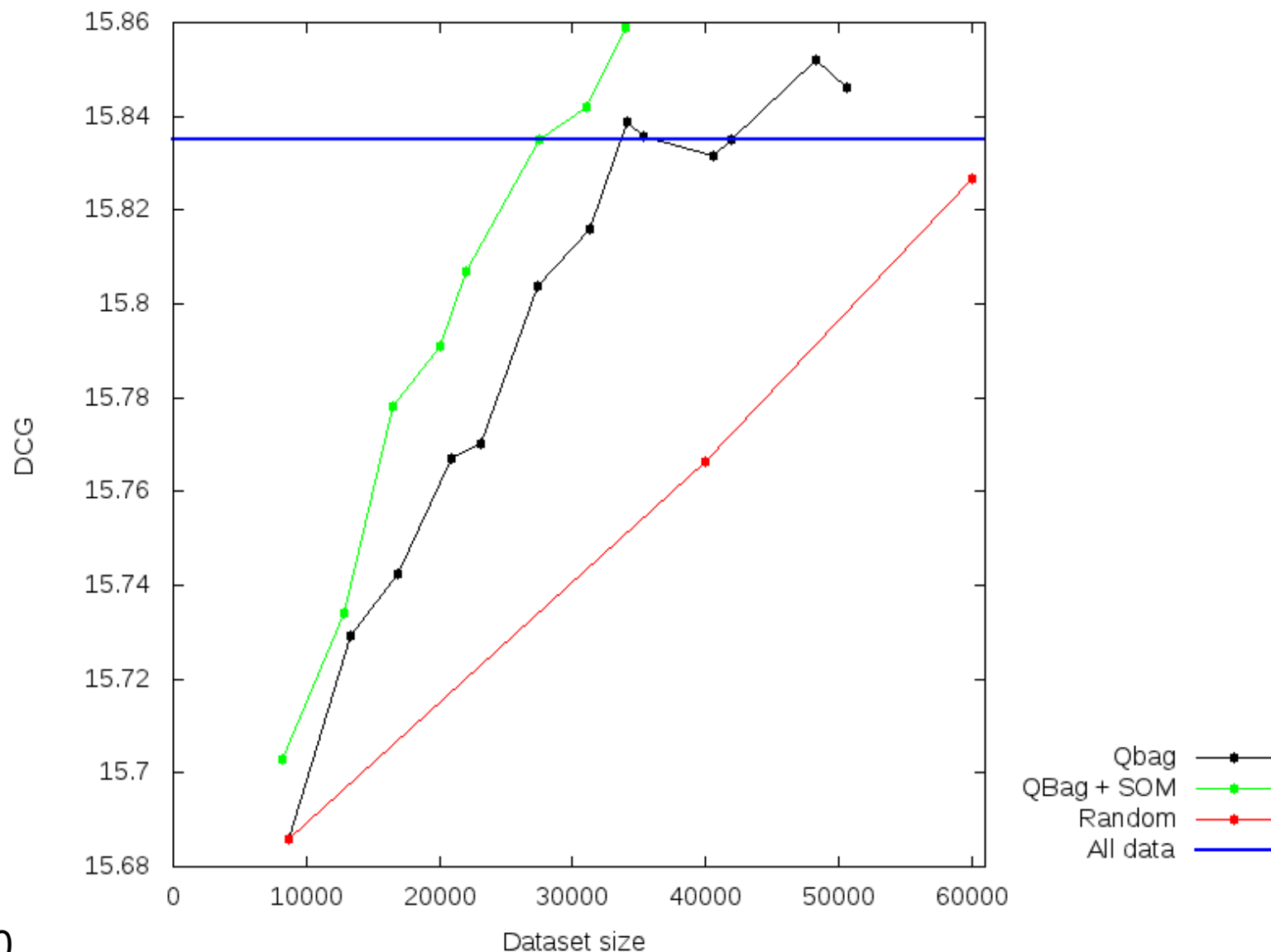
Example 2: Qbag + SOM heuristic

1. Transform regression to binary classification problem:
 - We should order pair of QDocuments
 - Apply *QBag* and select pairs of QDocuments
2. Apply SOM heuristic
 - Construct initial trainset T
 - Filter selected pairs



- A) Don't take pair
- B) May be
- C) Take pair

Example 2: Qbag + SOM heuristic: results

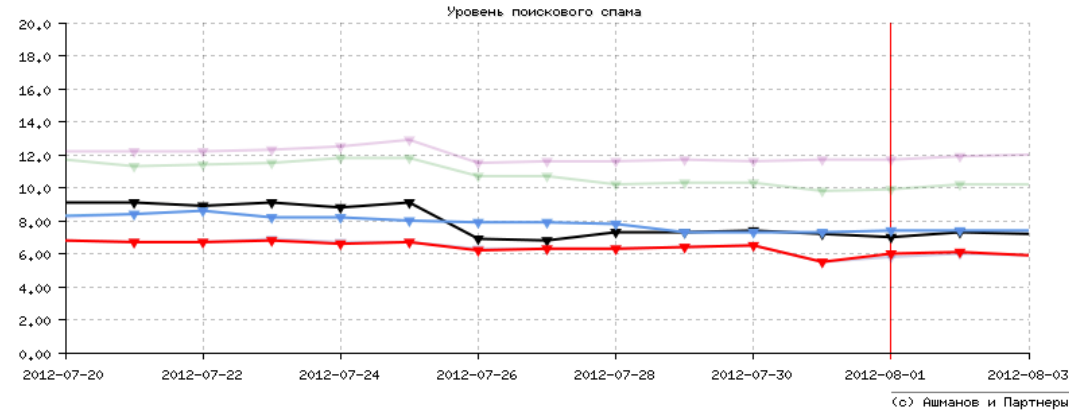


Example 3,4: Antispam, Antiporn



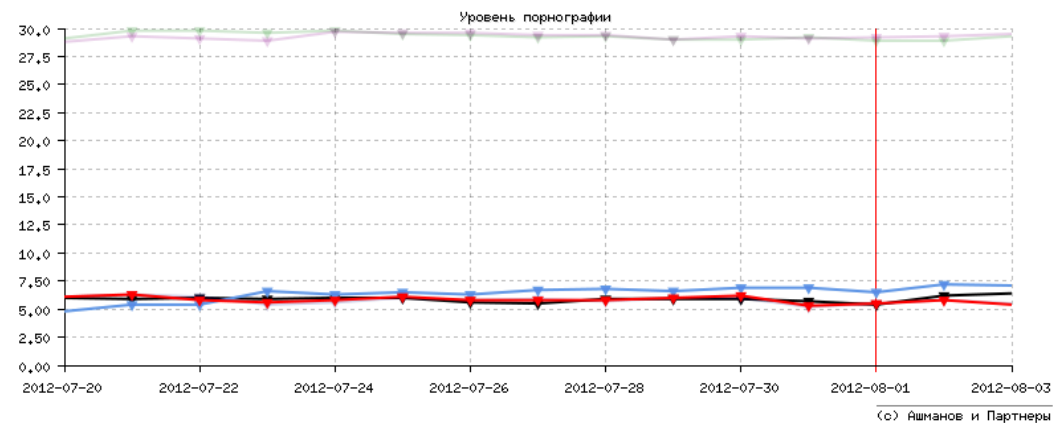
Antispam

- Neuron net
- SOM balancing



Antiporn

- Decision trees
- Uncertainty sampling



Thank you!

Reference:

- <http://active-learning.net/>
- Burr Settles. *Active Learning Literature Survey*. Computer Sciences Technical Report 1648, University of Wisconsin–Madison. 2009.