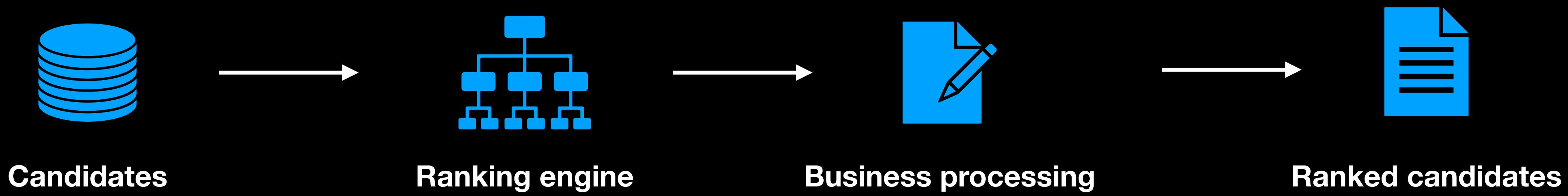


Recommendation Systems

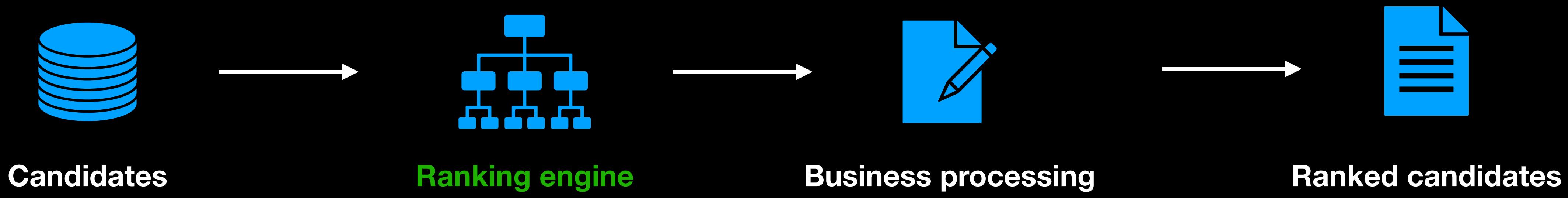
Unclassical embeddings + Learning to rank

Eugeniy Malyutin / Sergey Dudorov

RecSystem structure



RecSystem structure



Ok, ML

How to define it?

- ...
- $x \in X = \{x_1, \dots, x_n\}$ – items, $y \in Y = \{y_1, \dots, y_n\}$ – «labels»
- There are X_n, Y_n – given dataset with «answers», we believe that there is implicit dependency $y^* : X \rightarrow Y$
- Need to define algorithm $\alpha : X \rightarrow Y$.
- If $Y \in \{0,1\}$ – this task is called binary classification
If $Y \in \{y_0, y_1, \dots, y_n\}$ – this task is called multi class classification
If $Y = R$ – this task is called regression

Ok, ML

How to define it?

- ...
- $x \in X = \{x_1, \dots, x_n\}$ – items, $y \in Y = \{y_1, \dots, y_n\}$ – «labels»
- There are X_n, Y_n – given dataset with «answers», we believe that there is implicit dependency $y^* : X \rightarrow Y$
- Need to define algorithm $\alpha : X \rightarrow Y$.
- If $Y \in \{0,1\}$ – this task is called binary classification
If $Y \in \{y_0, y_1, \dots, y_n\}$ – this task is called multi class classification
If $Y = R$ – this task is called regression

Ok, binary classification.

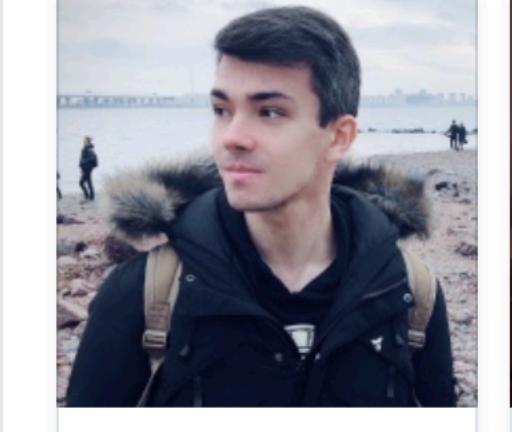
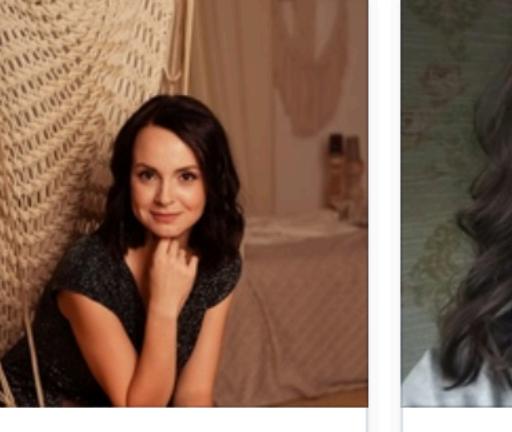
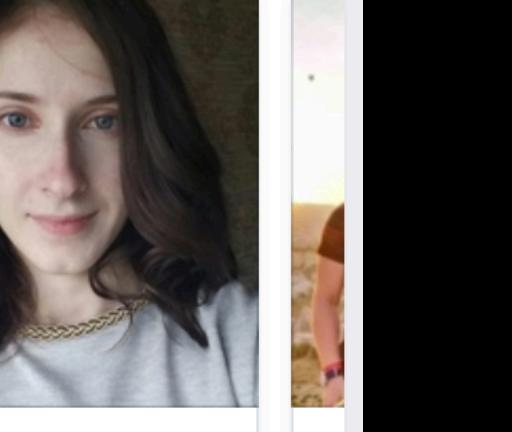
Does our cases fits binary classification?

Recommended

-  Wild Bags ✓
Open group
Follow
2 friends · 48,827 members
-  Пикабю ✓
Humor
Follow
58 friends · 3,057,449 followers
-  CoSSA
Mass media
Follow
12 friends · 165,595 followers
-  Ton Twitter
Humor
Follow
8 friends · 99,993 followers
-  PYE
Open group
Follow
6 friends · 42,068 members
-  Semrush
Websites
Follow
4 friends · 2,531 followers
-  Albertina
Gallery
Follow
19 friends · 950,042 followers

ВОЗМОЖНО, ВЫ ЗНАКОМЫ

Show all >

-  Vladislav Kaverin
ВКонтакте
29 общих друзей
+ Add
-  Natasha Belay
ВКонтакте
36 общих друзей
+ Add
-  Ksenia Timofeeva
Пермь
4 общих друга
+ Add

Customers who viewed this item also viewed

- 
Daily Ritual Women's
Lived-in Cotton Short-Sleeve Crewneck Maxi Dress
★★★★★ 42
\$23.40 - \$26.00
- 
Daily Ritual Women's
Jersey Mock-Neck Maxi Dress
★★★★★ 34
\$24.50
- 
Daily Ritual Women's
Jersey Sleeveless V-Neck Dress
★★★★★ 154
\$20.00
- 
Daily Ritual Women's
Jersey Crewneck Muscle Sleeve Maxi Dress with
Side Slit
★★★★★ 43
\$20.40 - \$24.50
- 
ZYX Women You are My
Sunshine Letter Print Tops
Casual Short Sleeve Tee
Rainbow T-Shirt
★★★★★ 1
\$16.98
- 
Daily Ritual Women's
Supersoft Terry Hooded
Short-Sleeve Sweatshirt
★★★★★ 16
\$25.20 - \$28.00
- 
Amazon Essentials
Women's Tank Maxi Dress
★★★★★ 24
\$26.00

Learning to rank

Definition

- X – set of objects. Exists: $x_i \prec x_j$ – **order** relation.

Our task to implement ranking function: $x_i \prec x_j \rightarrow a(x_i) < a(x_j)$

- Most of times order exists alongside with groupId (query relevant documents search engine)
if $y(q, d) = \{0,1\}$ – binary relevance, $d_q^{(i)}$ - i-th doc by $a(x_i)$ desc
- Ok, what's do you need more?
- Metrics
- Algorithms
- Tricks

Learning to rank

Metrics

- **Precision:** $P_n(q) = 1/n \sum_{i=1}^n y(q, d_q^{(i)})$ (percentage of correct in first n).
- Problems?
- **Average Precision:** $AP(q) = \sum y(d, d_q^{(n)})P_n(q)) / \sum y(q, d_n^q)$
(average P_n by all relevant docs position)
- Problems?
- **Mean average precision:** $MAP = 1/|Q| \sum_q AP(x)$

Learning to rank

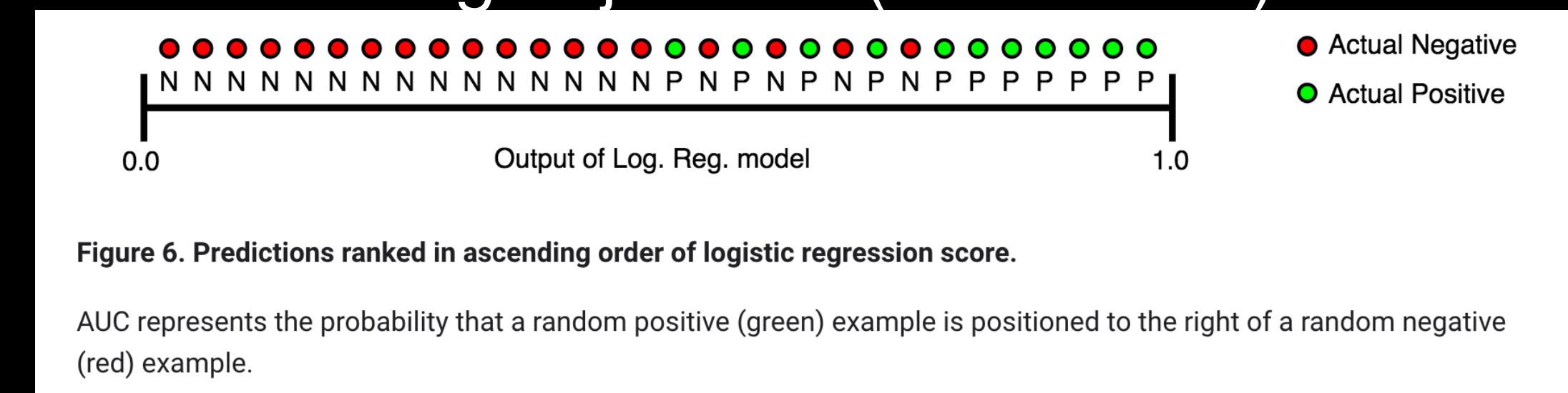
Metrics

- **MRR:** $MRR = 1/|Q| \sum_{i=1}^{|Q|} 1/rank_i$:
 - Sum over relevant!
 - Easy to interpret
- **DCG** (discounted cumulative gain): $DCG_n(q) = \sum_i^n G_q(d_q^i)D(i)$
 $G_q(d_q^i) = (y(d, q))$ – gain, 1/0 for relevant/unrelevant doc's
 $D(i) = 1/\log_2(i + 1)$ – discounts, higher relevant docs are – better
- **NDCG:** $NDCG(q) = DCG_n(q)/\max DCG_n(q)$ – normalized version

Learning to rank

Metrics

- **DCG** (discounted cumulative gain): $DCG_n(q) = \sum_i^n G_q(d_q^i)D(i)$
 $G_q(d_q^i) = (y(d, q))$ – gain, 1/0 for relevant/unrelevant doc's
 $D(i) = 1/\log_2(i + 1)$ – discounts, higher relevant docs are – better
- **NDCG:** $NDCG(q) = DCG_n(q)/\max DCG_n(q)$ – normalized version
- Also NDCG exists w.r.t. groupId, possible to use averaged over groups.
- Strongly correlates with AUC-per-groupId, but works with ranking objectives (unbounded)



Learning to rank

NDCG by example

NDCG - Example

4 documents: d_1, d_2, d_3, d_4

i	Ground Truth		Ranking Function ₁		Ranking Function ₂	
	Document Order	r_i	Document Order	r_i	Document Order	r_i
1	d_4	2	d_3	2	d_3	2
2	d_3	2	d_4	2	d_2	1
3	d_2	1	d_2	1	d_4	2
4	d_1	0	d_1	0	d_1	0
	$NDCG_{GT}=1.00$		$NDCG_{RF1}=1.00$		$NDCG_{RF2}=0.9203$	

$$DCG_{GT} = 2 + \left(\frac{2}{\log_2 2} + \frac{1}{\log_2 3} + \frac{0}{\log_2 4} \right) = 4.6309$$

$$DCG_{RF1} = 2 + \left(\frac{2}{\log_2 2} + \frac{1}{\log_2 3} + \frac{0}{\log_2 4} \right) = 4.6309$$

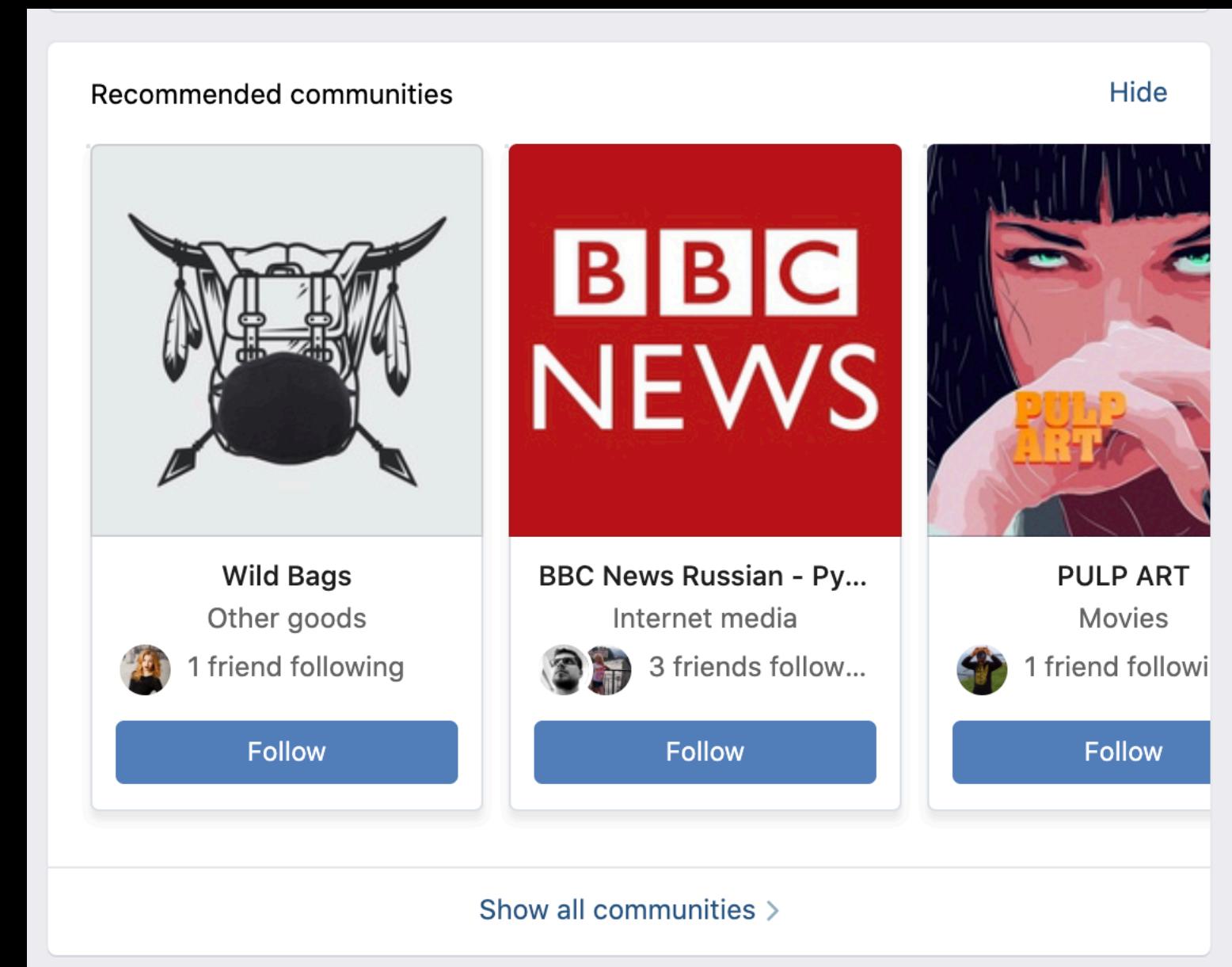
$$DCG_{RF2} = 2 + \left(\frac{1}{\log_2 2} + \frac{2}{\log_2 3} + \frac{0}{\log_2 4} \right) = 4.2619$$

$$MaxDCG = DCG_{GT} = 4.6309$$

Learning to rank

Metrics

- **DCG** (discounted cumulative gain): $DCG_n(q) = \sum_i^n G_q(d_q^i)D(i)$
 $G_q(d_q^i) = (y(d, q))$ – gain, 1/0 for relevant/unrelevant doc's
 $D(i) = 1/\log_2(i + 1)$ – discounts, higher relevant docs are – better
- **NDCG:** $NDCG(q) = DCG_n(q)/\max DCG_n(q)$ – normalized version
- Also popular **NDCG-at-k** (k=3, for example)



Learning to **rank**

Classification ontology

- **Pointwise:** sorting by score $\alpha(x_i)$ estimating «relevance»
- **Pairwise:** compares objects $a(x_i, x_j) :> 0 \text{ if } x_i < x_j$ — decide which one is more relevant
- **Listwise:** Takes **the entire list** of candidates and optimise it's order

Learning to rank

Classification ontology

- **Pointwise:** sorting by score $\alpha(x_i)$ estimating «relevance»

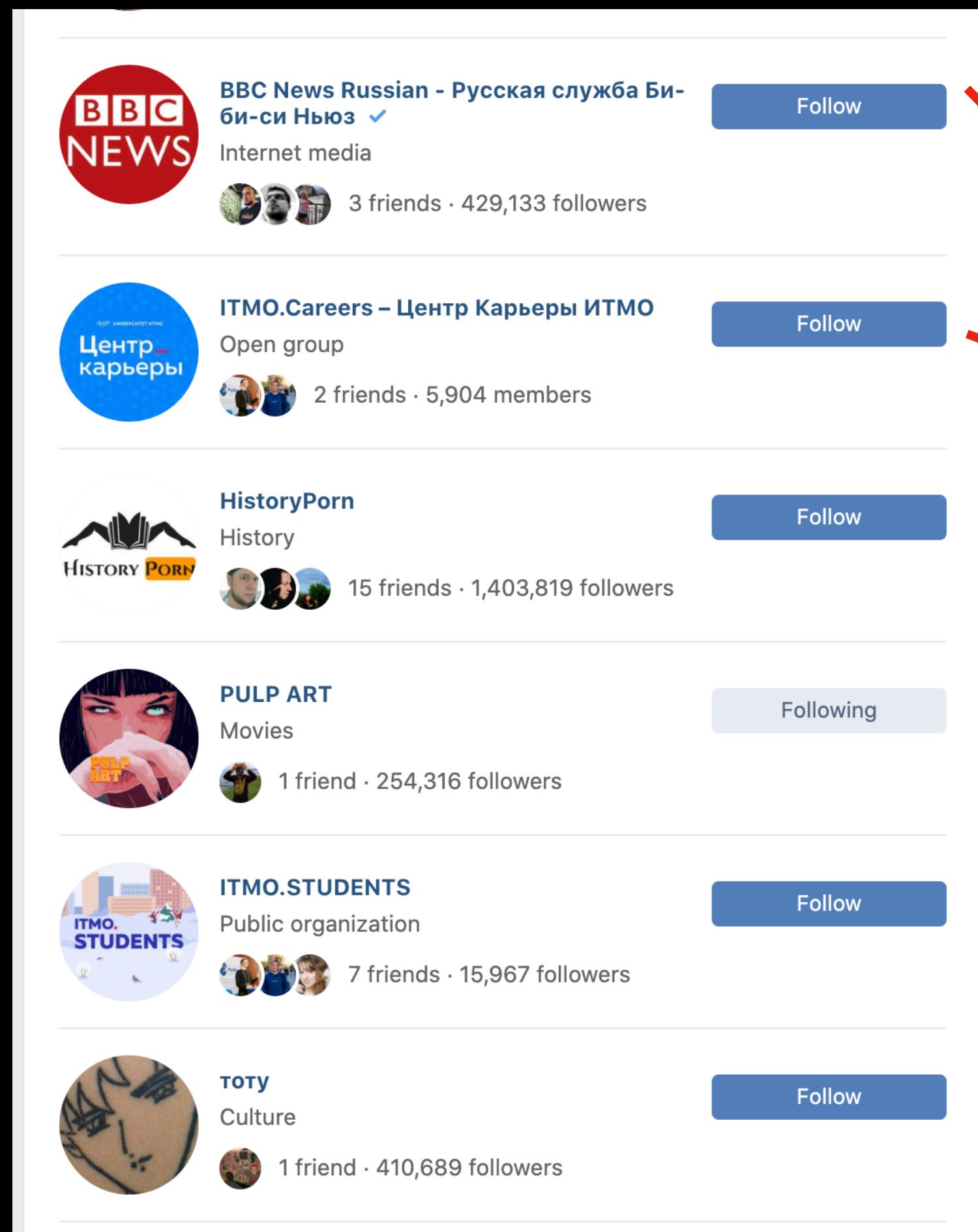


- **Positives** – $y_3 = 1$
- **Negatives** – $y_{0,1,2,4,5} = 0$
- $\alpha(x)$ gives relevance score $\in [0,1]$

Learning to rank

Classification ontology

- **Pairwise:** decide which one is more relevant



Negative

Negative

Negative

Positive

Negative

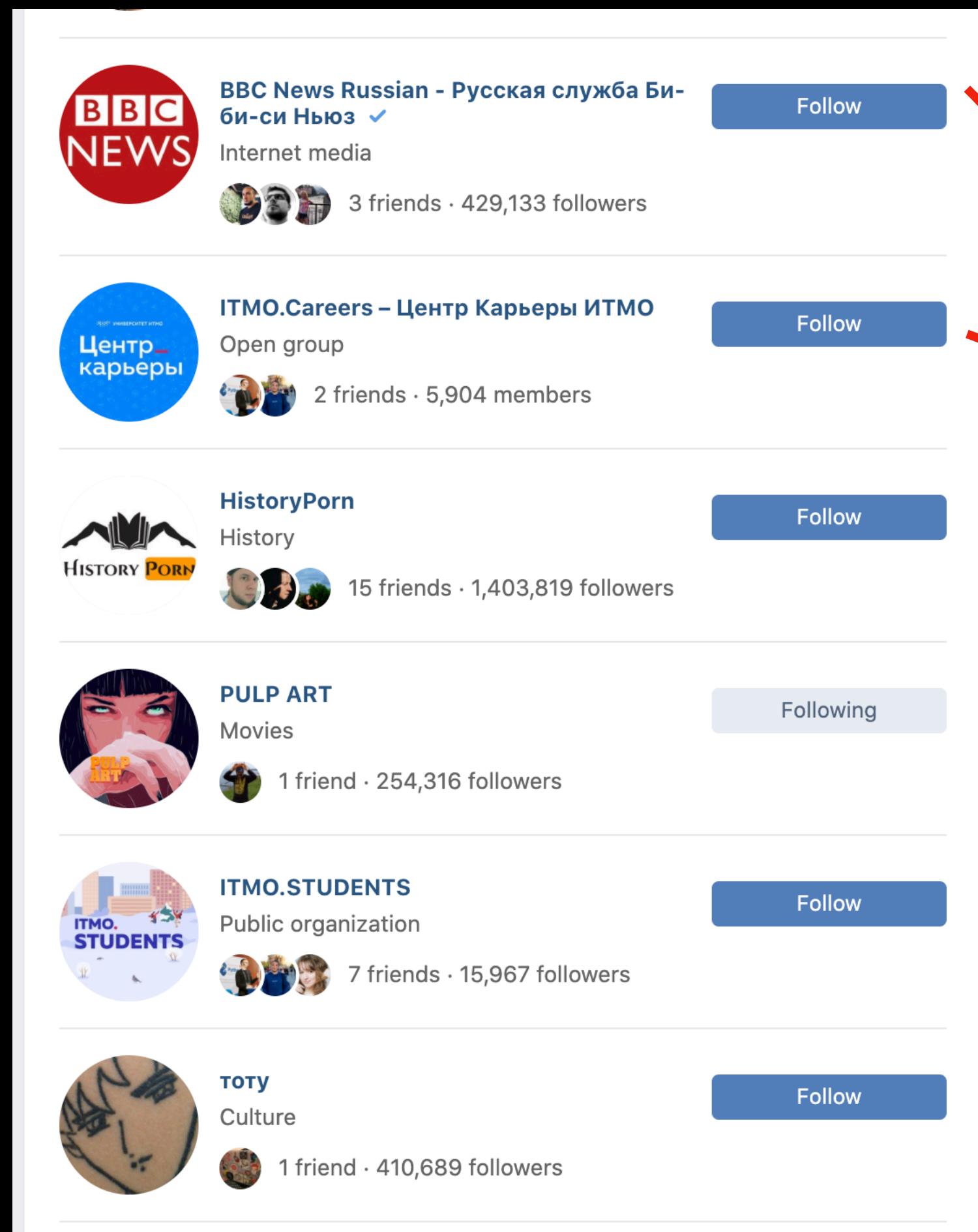
Negative

- **Train set:** $\{x_3 \prec x_0, x_3 \prec x_1, \dots, x_3 \prec x_5\}$
- $\alpha(x_i, x_j)$ compares pair of objects
- OK, what's a problem?

Learning to rank

Classification ontology

- **Pairwise:** decide which one is more relevant



Negative

Negative

Negative

Positive

Negative

Negative

- **Train set:** $\{x_3 \prec x_0, x_3 \prec x_1, \dots, x_3 \prec x_5\}$
- $\alpha(x_i, x_j)$ compares pair of objects
- OK, what's a problem?
- How much α applications do you need to sort 10000 - elements list?

Learning to rank

Classification ontology

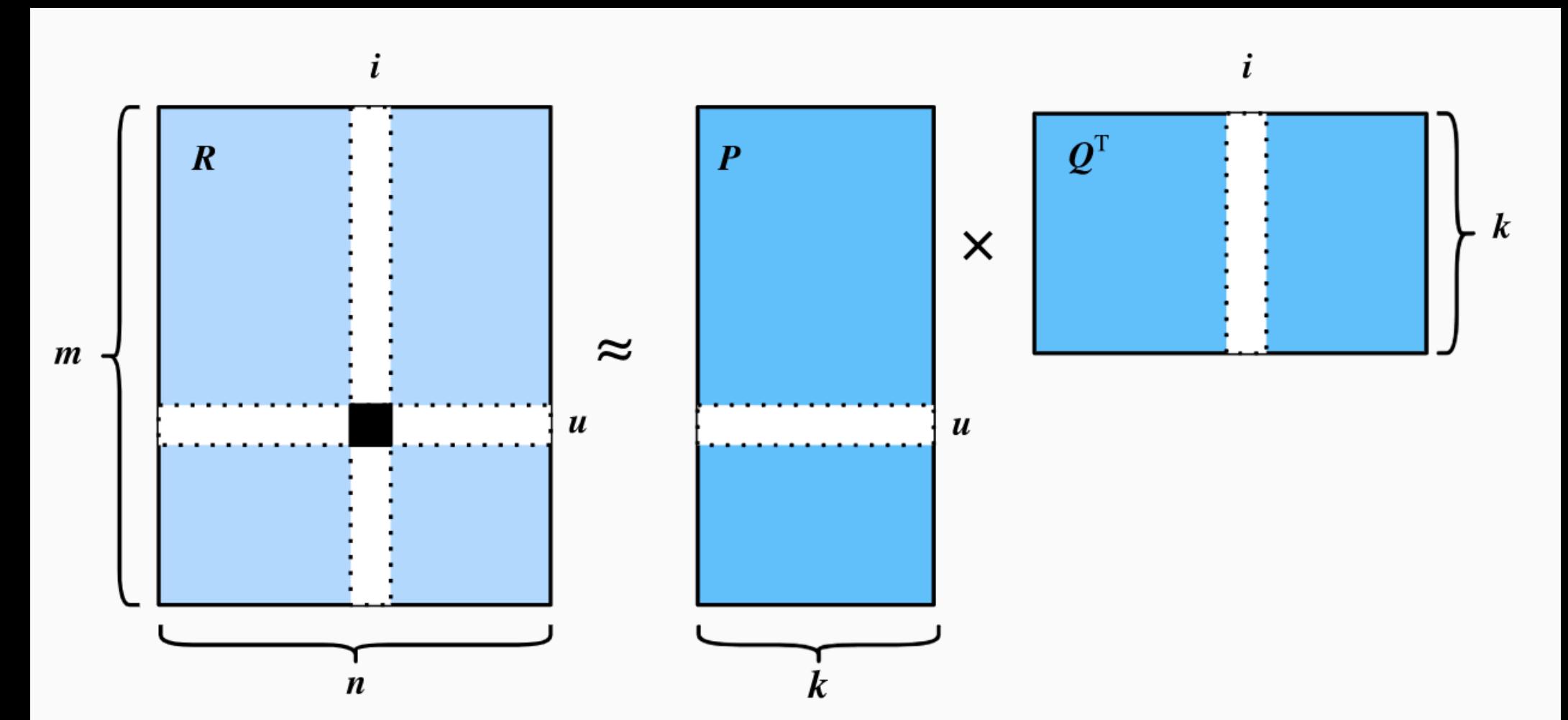
- **Listwise:** optimize the whole list



Learning to rank

SVD-factorisation

$$\operatorname{argmin}_{P,Q,b} \sum_{u,i \in K} ||R_{ui} - \hat{R}_{ui}||^2 + \lambda(||P||^2 + ||Q||^2 + b_u^2 + b_i^2)$$



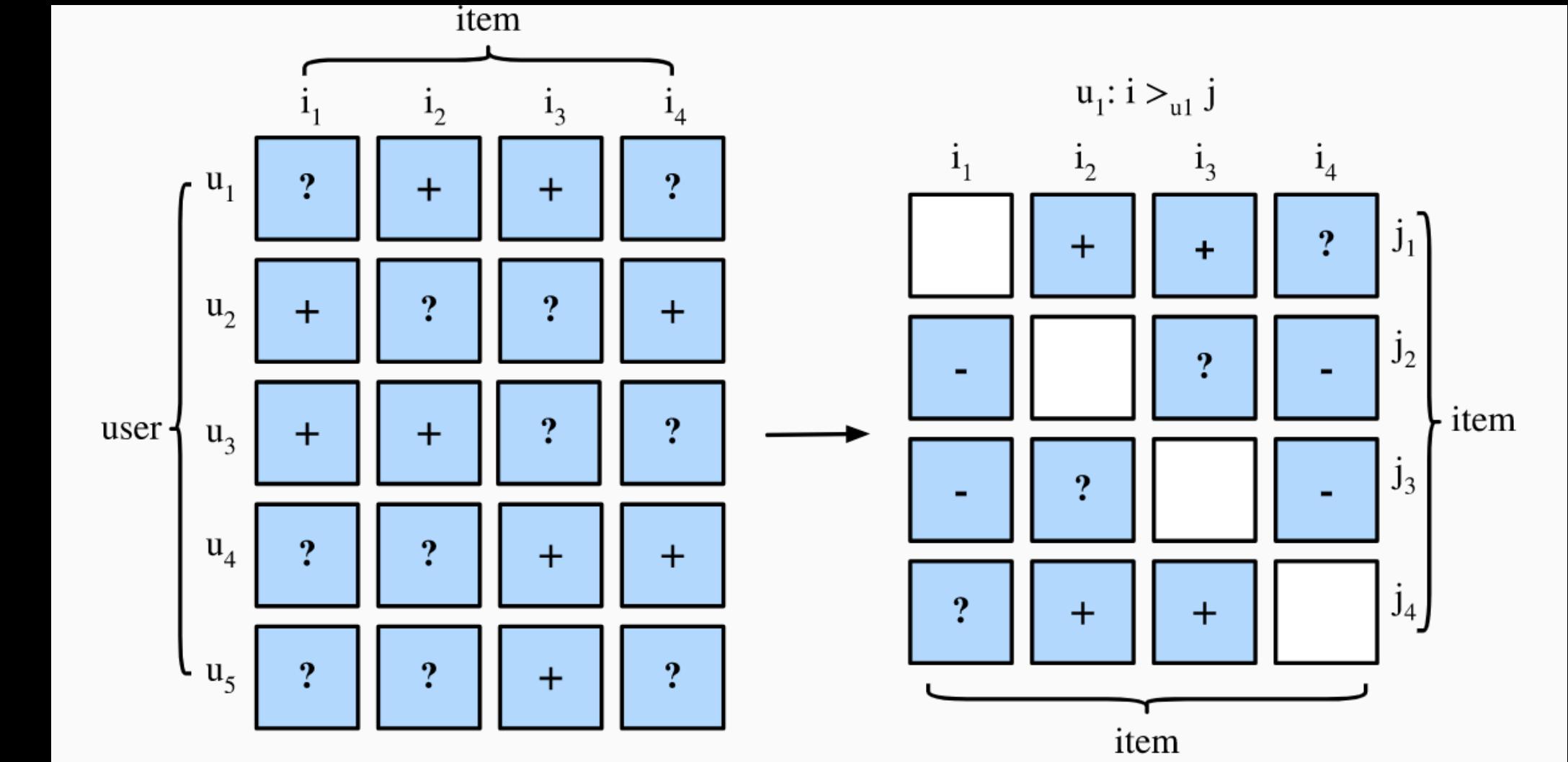
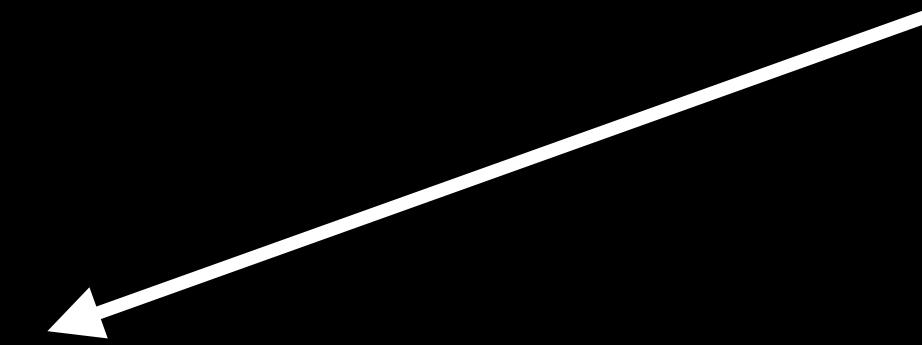
$$\operatorname{argmin}_{p,q,b} \sum_{u,i \in K} (\langle p_u, q_i \rangle + b_u + b_i + \mu - x_{ij})^2 + \alpha \sum_i ||u_i||^2 + \beta \sum_j ||p_j||^2 + \gamma ||b_u|| + \theta ||b_i||$$

$$\alpha(u, i) = \langle p_u, q_i \rangle + b_u + b_i + \mu \longrightarrow \text{Pointwise algorithm(!)}$$

Learning to rank

BPR-factorisation

$$p(\Theta | >_u) \propto p(>_u | \Theta)p(\Theta)$$



$$\text{BPR-OPT} := \ln p(\Theta | >_u)$$

$$\propto \ln p(>_u | \Theta)p(\Theta)$$

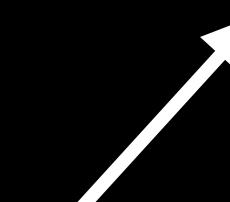
$$= \ln \prod_{(u,i,j \in D)} \sigma(\hat{y}_{ui} - \hat{y}_{uj})p(\Theta)$$

$$= \sum_{(u,i,j \in D)} \ln \sigma(\hat{y}_{ui} - \hat{y}_{uj}) + \ln p(\Theta)$$

$$= \sum_{(u,i,j \in D)} \ln \sigma(\hat{y}_{ui} - \hat{y}_{uj}) - \lambda_\Theta \|\Theta\|^2$$

$$r_{ui} = \langle e_u, e_i \rangle + b_u + b_i + \mu$$

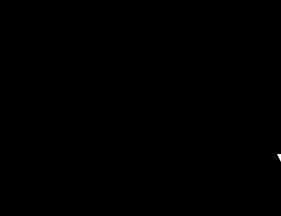
$$\alpha(x_i, x_j, u) = \ln(\sigma(\langle e_u, e_i \rangle + b_u + b_i - \langle e_u, e_j \rangle - b_u - b_j))$$



$$\alpha(u, i) = \langle e_u, e_i \rangle$$

Sort with α

PAIRwise algorithm while learn
POINTwise algorithm on runtime!



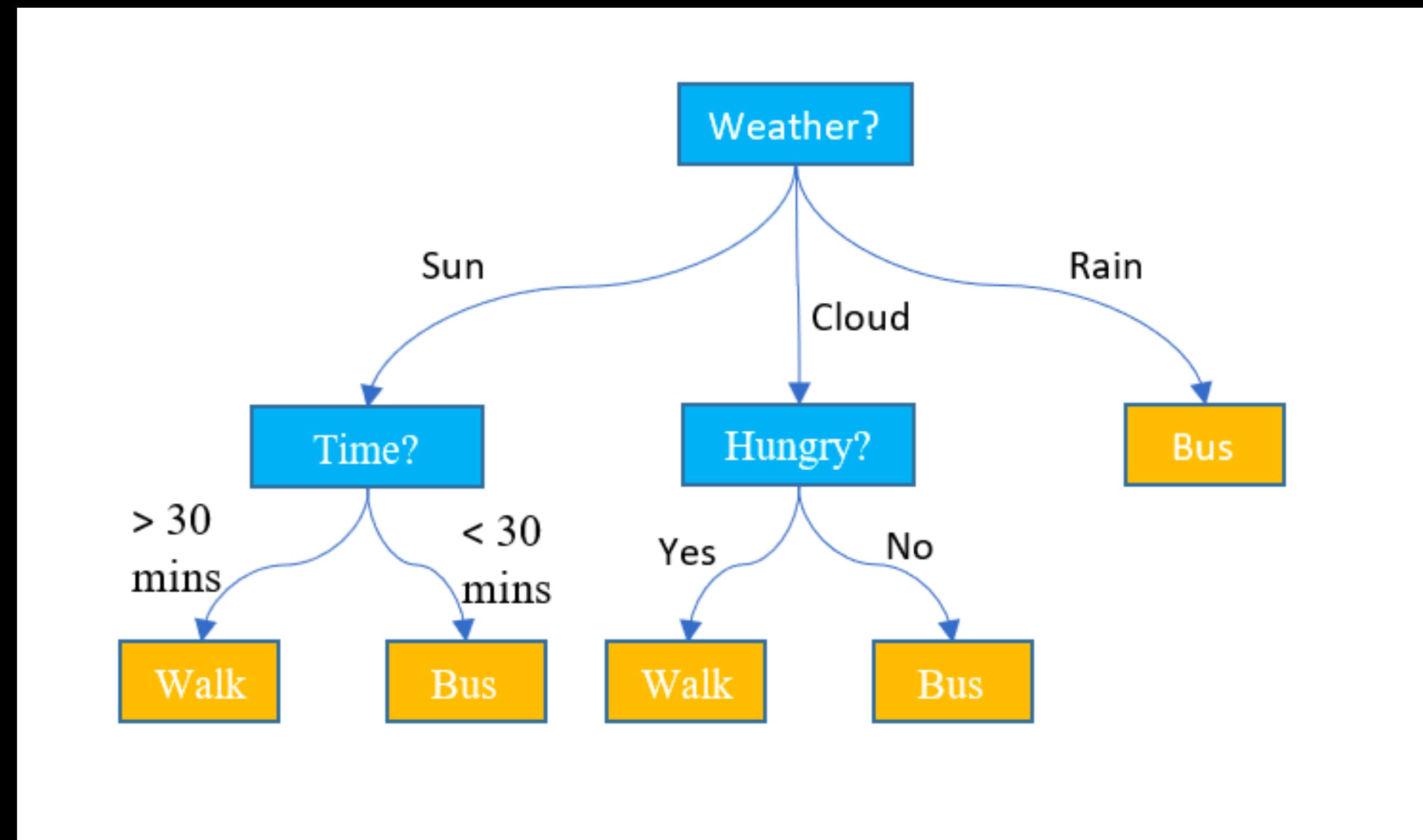
Learning to rank

Summary

- Ranking **differs** either classification or regression due to strict **order properties** on **relevance**
- **It** leads us to new metrics: MAP, MRR, AUC(\sim), NDCG
- Algorithms ontology: **point-**, **pair-**, **list**-wise
- There are **no** working «fair» **listwise** algorithms **IRL** (at least author saw no one)
- **Best pairwise** algorithms – camouflages themselves as **pointwise**.

Ranking algorithms

Tree reminder



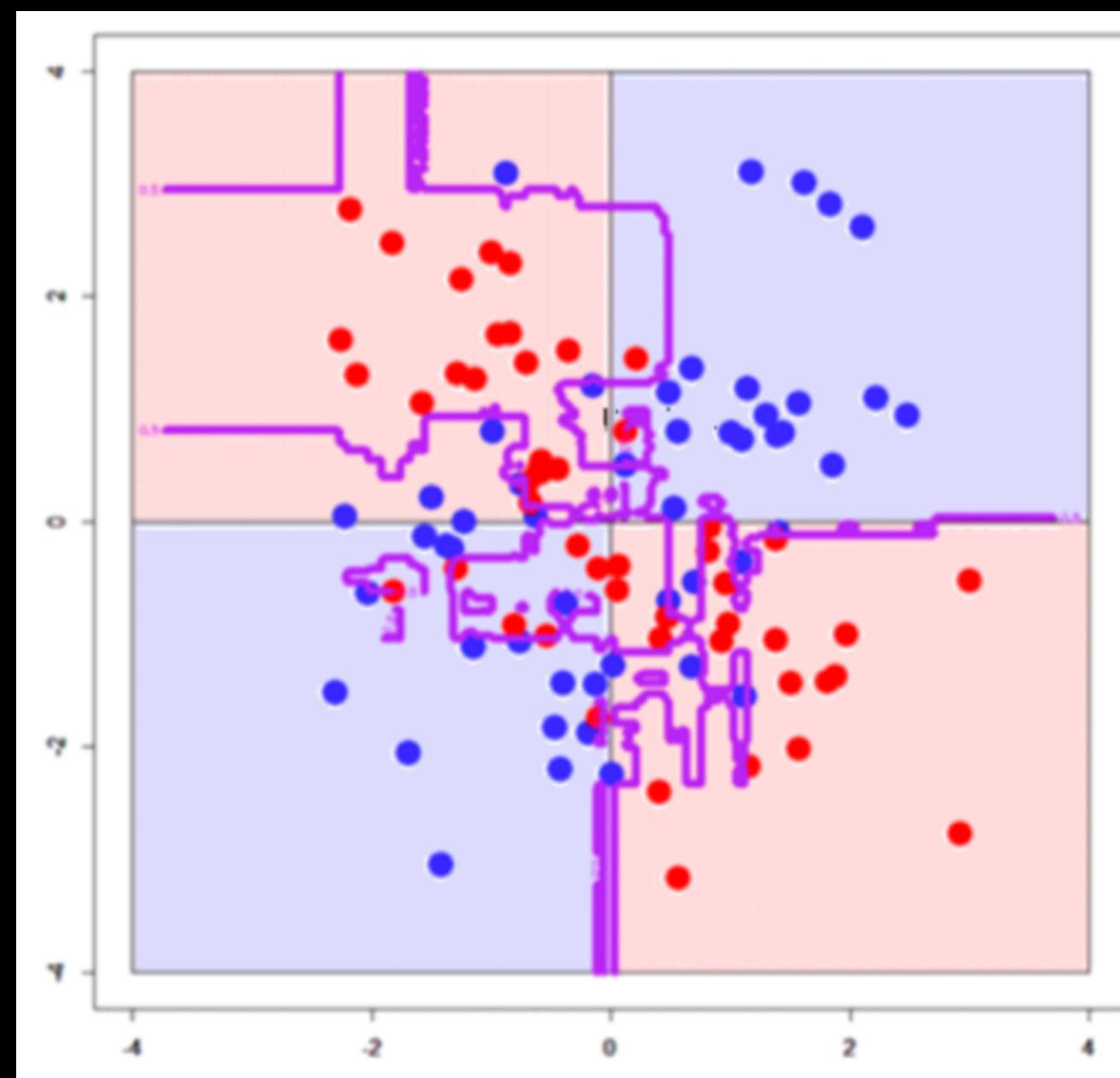
Ranking algorithms

Splitting criteria:

- **Split:** $[x_j < t]$
- **Loss criteria:** $Q(X_m, j, t) = \frac{|X_l|}{X_M} H(X_l) + \frac{|X_r|}{|X_m|} H(X_m)$
- **Information criteria:**
 - **Regression:** $H(x) = 1/|X_m| \sum (y_i - \bar{y}(X))^2$
 - **Gini:** $p_k = \frac{1}{|X|} \sum_{i \in X} [y_i = k]$ and $H(x) = \sum_{k=1}^K p_k(1 - p_k)$
 - **Entropy:** $H(X) = \sum_{k=1}^K p_k \ln(p_k)$

Ranking algorithms

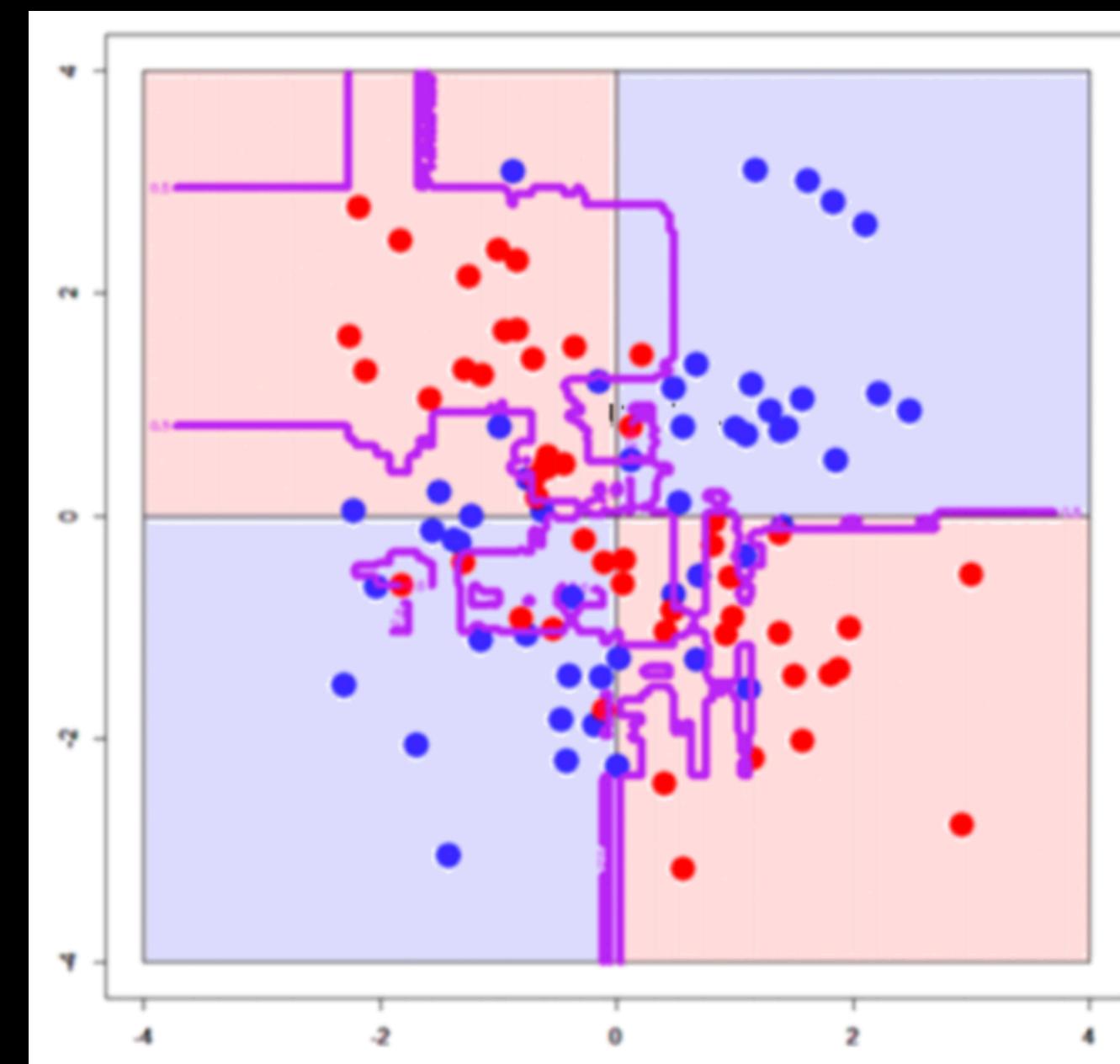
What's went wrong?



Ranking algorithms

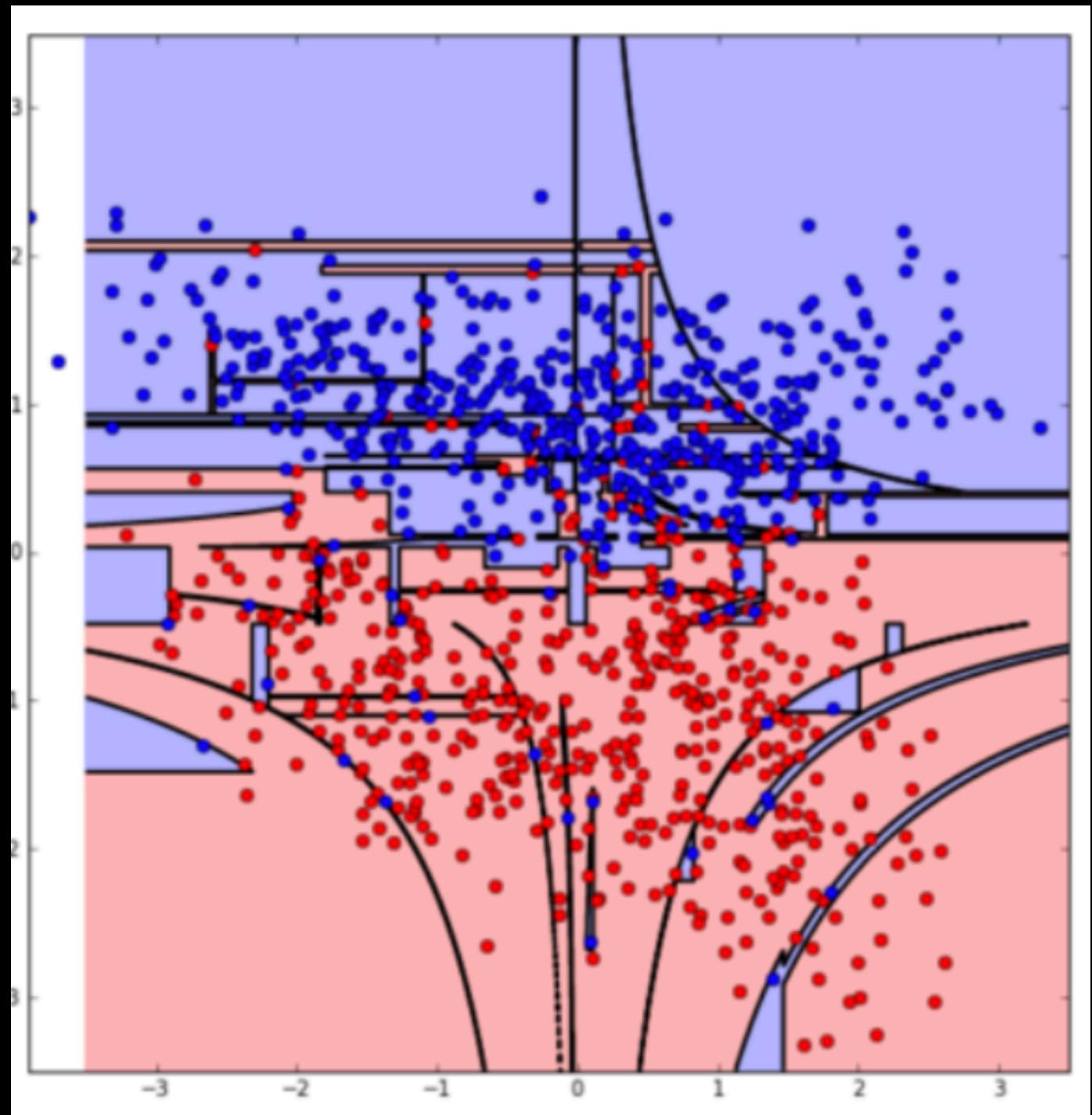
What's went wrong?

- Stop when:
 - Max length
 - Min elements
 - Pruning

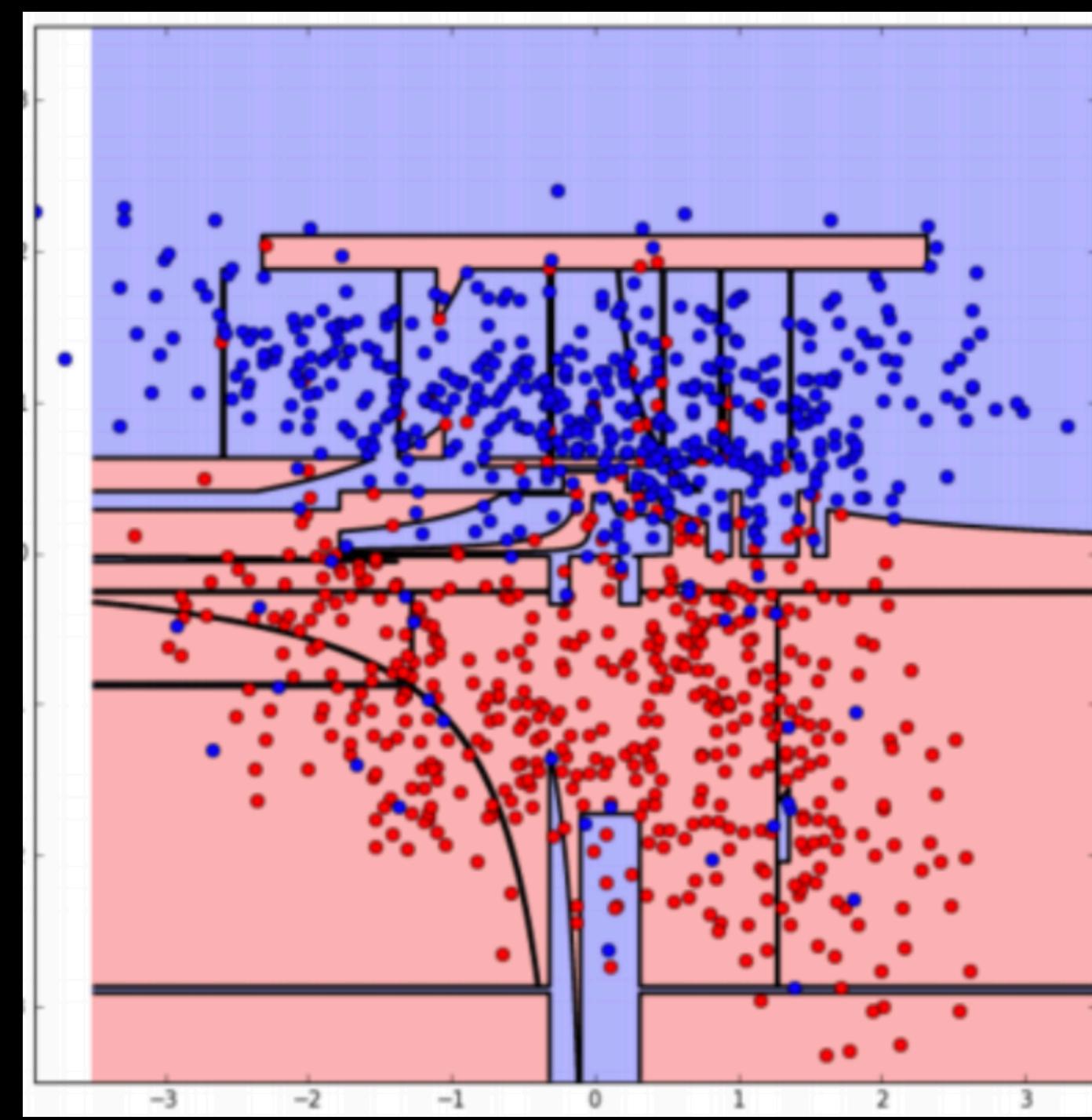


Ranking algorithms

What's went wrong?



Overfitted



Overfitted and unstable

Ranking algorithms

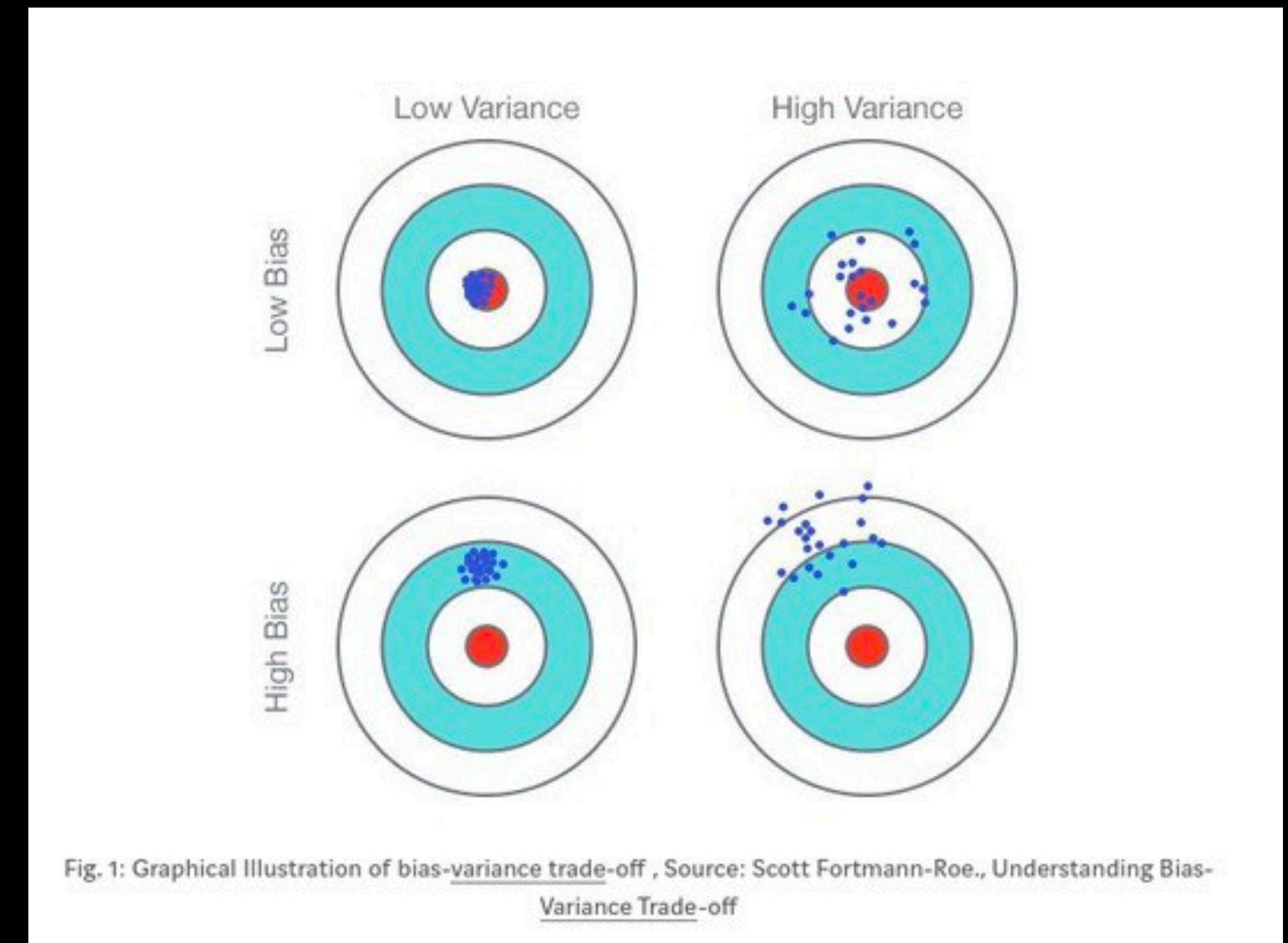
Error decomposition

- What is?

Ranking algorithms

Error decomposition

- **Error:**
 - **Noise:** world's imperfectness measure, exists even for ideal model on ideal data
 - **Bias:** deviation of concrete model from ideal one (averaged over all data sets)
 - **Variance:** dispersion of models answers caused by different datasets



Ranking algorithms

Composition idea

- Takes N different trees
- **Averages** answers:
 - **Regression:** $a(x) = \sum_1^N b_n(x)$
 - **Classification:** most popular answer
 - **PRBLMS?**

Ranking algorithms

Error decomposition

- **Trees:**
 - High variance
 - Low bias
- **Linear algorithms:**
 - Low variance
 - High bias
- **Composition:**
 - Same bias
 - $\text{Variance} = \frac{1}{N} (\text{algorithm_variance}) + \text{correlation}$
 - Decrease variance **N times** in case of **independent algorithms** — that's why **we need randomisation**

Ranking algorithms

Randomisation

- **Bagging (bootstrap)**: learn algorithms on random subsamples from train set. Less subsamples — more randomised trees.
- **Random sub-spaces**: random sub-set of features for each tree
- **Extreme randomised trees**: choose random subset of features on each split

Ranking algorithms

Summary

- **Pros:**
 - Powerful
 - Easy to parallel
 - With bootstrap 1/3 of data samples (for concrete tree) wasn't in train set — we can estimate test metrics on them (Out-of-bag score)
- **Cons:**
 - More trees -> more computational resources
 - Undirected search

Ranking algorithms

Boosting idea

- $b_0(x)$ – initial algorithm (zero, mean class, mean value)
- $a_m(x) = \sum_{i=1}^m b_i(x)$ – step of composition
- $F = \sum_{i=1}^N L(y_i, a_{m-1}(x_i) + b_i(x_i)) \rightarrow \min_b$ – optimisation
- $s = (s_1, s_2 \dots s_N)$ – displacements vector then
$$F = \sum_{i=1}^N L(y_i, a_{m-1}(x_i) + s_i) \rightarrow \min_s$$
- Optimal shift – $\nabla F = [\frac{dF}{da_{m-1}}(x_i)]_{i=1}^N = [\frac{\sum_{i=1}^N L(y_i, a_{m-1}(x_i))}{da_{m-1}}]_{i=1}^N = [\frac{dL(y_i, a_{m-1})}{da_{m-1}} x_i]_{i=1}^N$

Ranking algorithms

Boosting idea

- $a_m(x) = \sum_{i=1}^m b_i(x)$ – step of composition, $F = \sum_{i=1}^N L(y_i, a_{m-1}(x_i) + b_i(x_i)) \rightarrow \min_b$ – optimisation
- Optimal shift
$$s_i = -\nabla F = \left[\frac{dF}{da_{m-1}}(x_i) \right]_{i=1}^N = \left[\frac{\sum_{i=1}^N L(y_i, a_{m-1}(x_i))}{da_{m-1}}(x_i) \right]_{i=1}^N = \left[\frac{dL(y_i, a_{m-1})}{da_{m-1}}(x_i) \right]_{i=1}^N$$
- But wait, s_i is not an algorithm, it's a vector of numbers!
- Yep, but we can learn out algorithm $b_m(x_i) \rightarrow s_i$

Ranking algorithms

Summary

- Powerful (extremely powerful)
- Allows all this tricks of gradient algorithms: learning rate, decay, ...
- Hard to interpret (~)
- Works mostly with weak algorithms
- Allows you to set various range of functions(?) as loss

XGboost

XGBoost

$$\hat{y}_i^{(0)} = 0$$

$$\hat{y}_i^{(1)} = f_1(x_i) = \hat{y}_i^{(0)} + f_1(x_i)$$

$$\hat{y}_i^{(2)} = f_1(x_i) + f_2(x_i) = \hat{y}_i^{(1)} + f_2(x_i)$$

⋮

$$\hat{y}_i^{(t)} = \sum_{k=1}^t f_k(x_i) = \hat{y}_i^{(t-1)} + f_t(x_i)$$



$$\begin{aligned}\text{obj}^{(t)} &= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i) \\ &= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \text{constant}\end{aligned}$$

Composition structure

Objective function

$\Omega(f_i)$ – regularisation

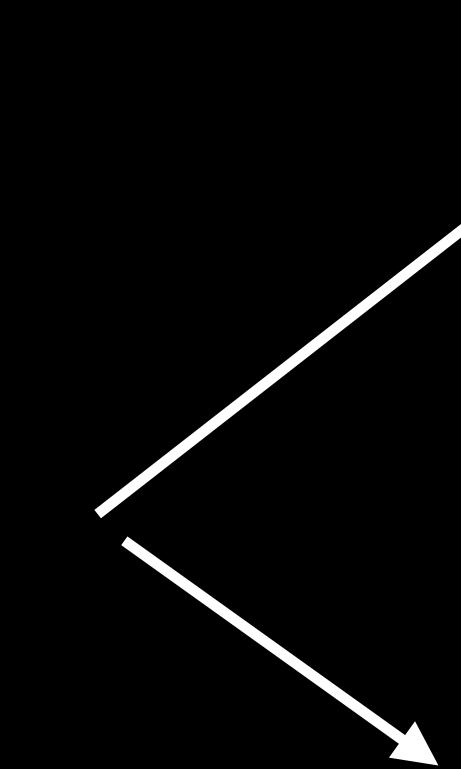
XGboost

XGBoost

$$\text{obj}^{(t)} = \sum_{i=1}^n l(y_i, \hat{y}_i^{(t)}) + \sum_{i=1}^t \Omega(f_i)$$

$$= \sum_{i=1}^n l(y_i, \hat{y}_i^{(t-1)} + f_t(x_i)) + \Omega(f_t) + \text{constant}$$

Objective function



$$\begin{aligned}\text{obj}^{(t)} &= \sum_{i=1}^n (y_i - (\hat{y}_i^{(t-1)} + f_t(x_i)))^2 + \sum_{i=1}^t \Omega(f_i) \\ &= \sum_{i=1}^n [2(\hat{y}_i^{(t-1)} - y_i)f_t(x_i) + f_t(x_i)^2] + \Omega(f_t) + \text{constant}\end{aligned}$$

MSE example

$$\text{obj}^{(t)} = \sum_{i=1}^n [l(y_i, \hat{y}_i^{(t-1)}) + g_i f_t(x_i) + \frac{1}{2} h_i f_t^2(x_i)] + \Omega(f_t) + \text{constant}$$

$$g_i = \partial_{\hat{y}_i^{(t-1)}} l(y_i, \hat{y}_i^{(t-1)})$$

$$h_i = \partial_{\hat{y}_i^{(t-1)}}^2 l(y_i, \hat{y}_i^{(t-1)})$$

2nd order Taylor's approximation

XGboost

Regularisation

$$f_t(x) = w_{q(x)}, w \in R^T, q : R^d \rightarrow \{1, 2, \dots, T\}.$$

Formal tree

- w – vector of scores on j -th index
- q – function to assign object to leaf
- T – number of leaves

$$\Omega(f) = \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

Complexity

XGboost

Regularisation in learning

$$\begin{aligned}\text{obj}^{(t)} &\approx \sum_{i=1}^n [g_i w_{q(x_i)} + \frac{1}{2} h_i w_{q(x_i)}^2] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma T\end{aligned}$$



$$\begin{aligned}\text{obj}^{(t)} &= \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T \\ G_j &= \sum_{i \in I_j} g_i & H_j &= \sum_{i \in I_j} h_i\end{aligned}$$

Obj function with Taylor and formal tree

«Simplification»

XGboost

...

$$\begin{aligned}\text{obj}^{(t)} &\approx \sum_{i=1}^n [g_i w_{q(x_i)} + \frac{1}{2} h_i w_{q(x_i)}^2] + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\ &= \sum_{j=1}^T [(\sum_{i \in I_j} g_i) w_j + \frac{1}{2} (\sum_{i \in I_j} h_i + \lambda) w_j^2] + \gamma T\end{aligned}$$



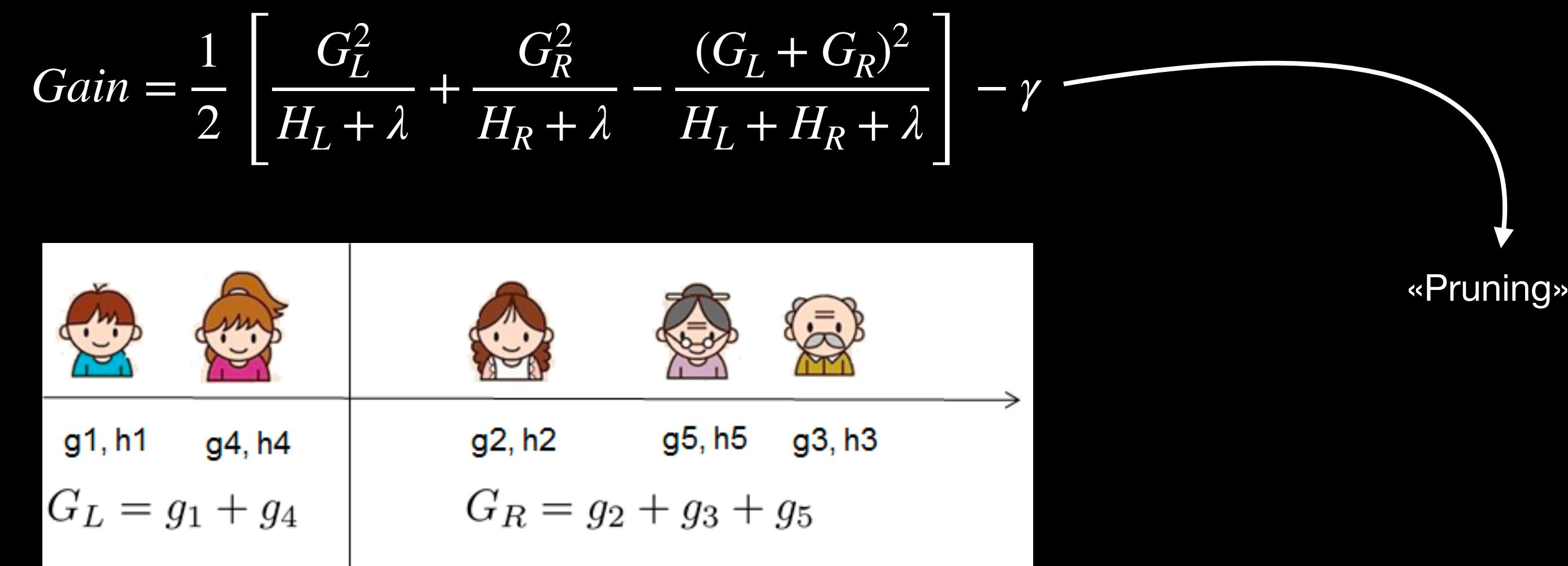
$$\begin{aligned}\text{obj}^{(t)} &= \sum_{j=1}^T [G_j w_j + \frac{1}{2} (H_j + \lambda) w_j^2] + \gamma T \\ G_j &= \sum_{i \in I_j} g_i & H_j &= \sum_{i \in I_j} h_i\end{aligned}$$

Obj function with Taylor and formal tree

«Simplification»

XGboost:

Sorting schema:



XGboost

Summing up:

- Decomposes loss-function in Taylor row
- Allows arbitrary set of loss functions
- Zip regularisation inside learning process
- Scan all dataset ones to find best split

XGboost

And his friends:

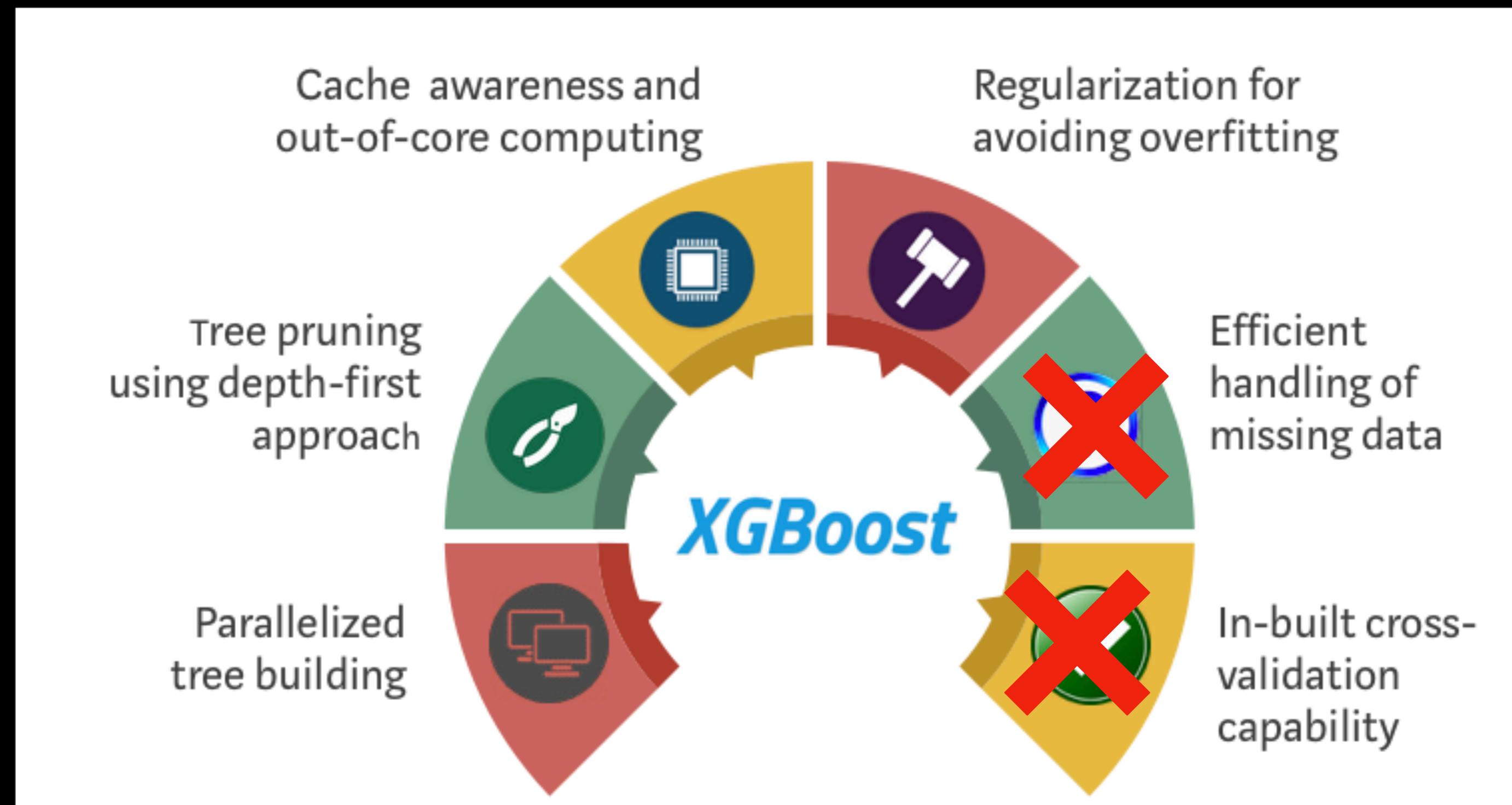
Pros:

- Well-known (est. 2016)
- Available on ton different platforms (as cpp binding, mostly)
- Allows distribute training (spark, for example)

Cons:

- Not so fancy comparing to IGBM, catBoost
- Not so powerful (~)
- Very strange sparsity

XGBoost:



XGBoost notes:

- XGBoost internal feature significance sucks, **use SHAP values**.
- XGBoost objectives:
 - **binary:logistic**: logistic regression for binary classification, output probability
 - **reg:squarederror**: regression with squared loss. Practically better on classification oO
 - **rank:{pairwise/map/ndcg}**: lambda mart ranking. Sometimes works significantly better than regr losses.
- Hyperparameters matters, tune them wisely (look refs).
- High gamma -> feature selection -> low gamma and features removed

Learning to rank

«Lambda»-smth

$$w = w + \eta \frac{\sigma}{1 + \exp(\sigma \langle x_i - x_j, w \rangle)} (x_i - x_j)$$

Gradient step optimising logit function on linear algorithm

$$w = w + \eta \frac{\sigma}{1 + \exp(\sigma \langle x_i - x_j, w \rangle)} |\Delta NDCG_{ij}| (x_i - x_j)$$

NDCG delta
by changing x_i and x_j

Gradient step optimising NDCG function on linear algorithm

Learning to rank

Summary

- Learning to rank - specific task with specific metrics and tricks.
- There are: point- pair- and list-wise algorithms.
- Best pair-wise algorithms looks as pointwise at runtime.
- There are tons ranking algorithms. Most powerful and/or used – GBDT
- XGBoost – allows arbitrary set of loss functions, incorporates regularisation into training, use Taylor 2nd order to estimate loss function.
- It's possible to optimise ranking metrics directly using lambda- techniques

Vector embeddings

Eugeny Malyutin



KING

QUEEN



QUEEN

Previously:

motel [0 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0] AND
hotel [0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 0] = 0

Поправь doc frequency Малютин

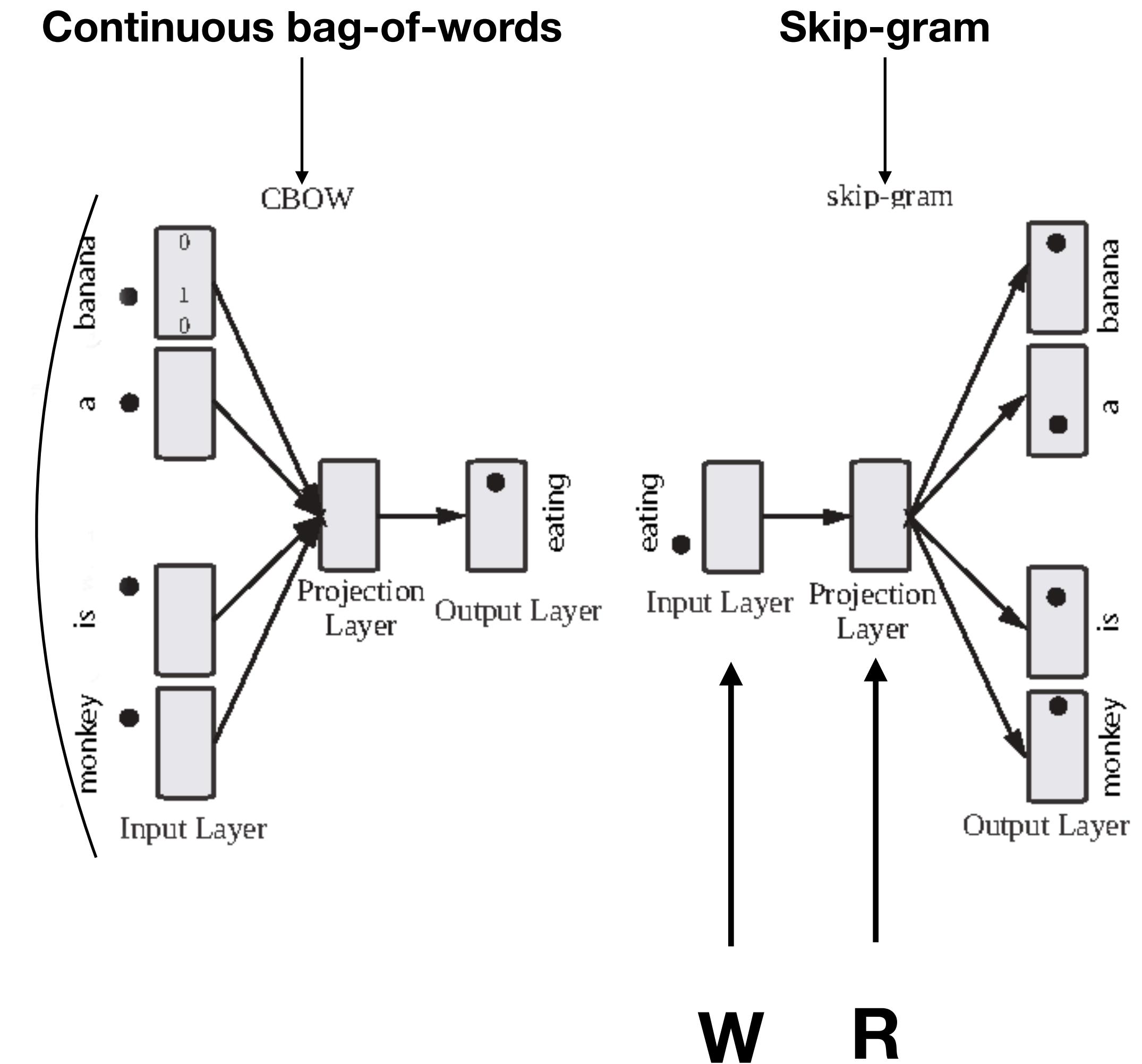
$$TF - IDF(w, d, C) = \frac{count(w, d)}{count(d)} * \log\left(\frac{\sum_{d' \in C} 1(w, d')}{|C|}\right)$$

One-hot encoding drawback:

- “monkeys eat bananas” and “apes consume fruits” - similarity equals to 0
- «Pouteria is widespread throughout the tropical regions of the world and monkeys eat their fruits»(c). What is Pouteria? Is it a tree?
- «a word is characterized by the company it keeps» – John Rupert Firth
- Ideally, we want vector representations where ***similar words*** end up with ***similar vectors***. **Dense** vectors. And when I say similar a mention some **similarity measure** (cosine).
- Even better, we’d want more similar representations when the words share some properties such as if they’re both plural or singular, verbs or adjectives or if they both reference to a male.

Word2vec scheme:

- It has an input layer that receives \mathbf{D} one-hot encoded words which are of dimension \mathbf{V} (the size of the vocabulary).
- It «averages» them, creating a single input vector.
- That input vector is multiplied by a weights matrix \mathbf{W} (that has size $\mathbf{V} \times \mathbf{D}$, being \mathbf{D} nothing less than the dimension of the vectors that you want to create). That gives you as a result a \mathbf{D} -dimensional vector.
- The vector is then multiplied by another matrix (\mathbf{R} - reverse \mathbf{W}), this one of size $\mathbf{D} \times \mathbf{V}$. The result will be a new \mathbf{V} -dimensional vector.
- That \mathbf{V} -dimensional vector is normalized to make all the entries a number between 0 and 1, and that all of them sum 1, using the **softmax function**, and that's the output. It has in the i -th position the predicted probability of the i -th word in the vocabulary of being the one in the middle for the given context.



The skipgram model

- We assume that, given the central target word, the context words are generated independently of each other.

$$P(\text{the, man, his, son} \mid \text{loves}) = P(\text{the} \mid \text{loves}) * P(\text{man} \mid \text{loves}) * P(\text{his} \mid \text{loves}) * P(\text{son} \mid \text{loves})$$

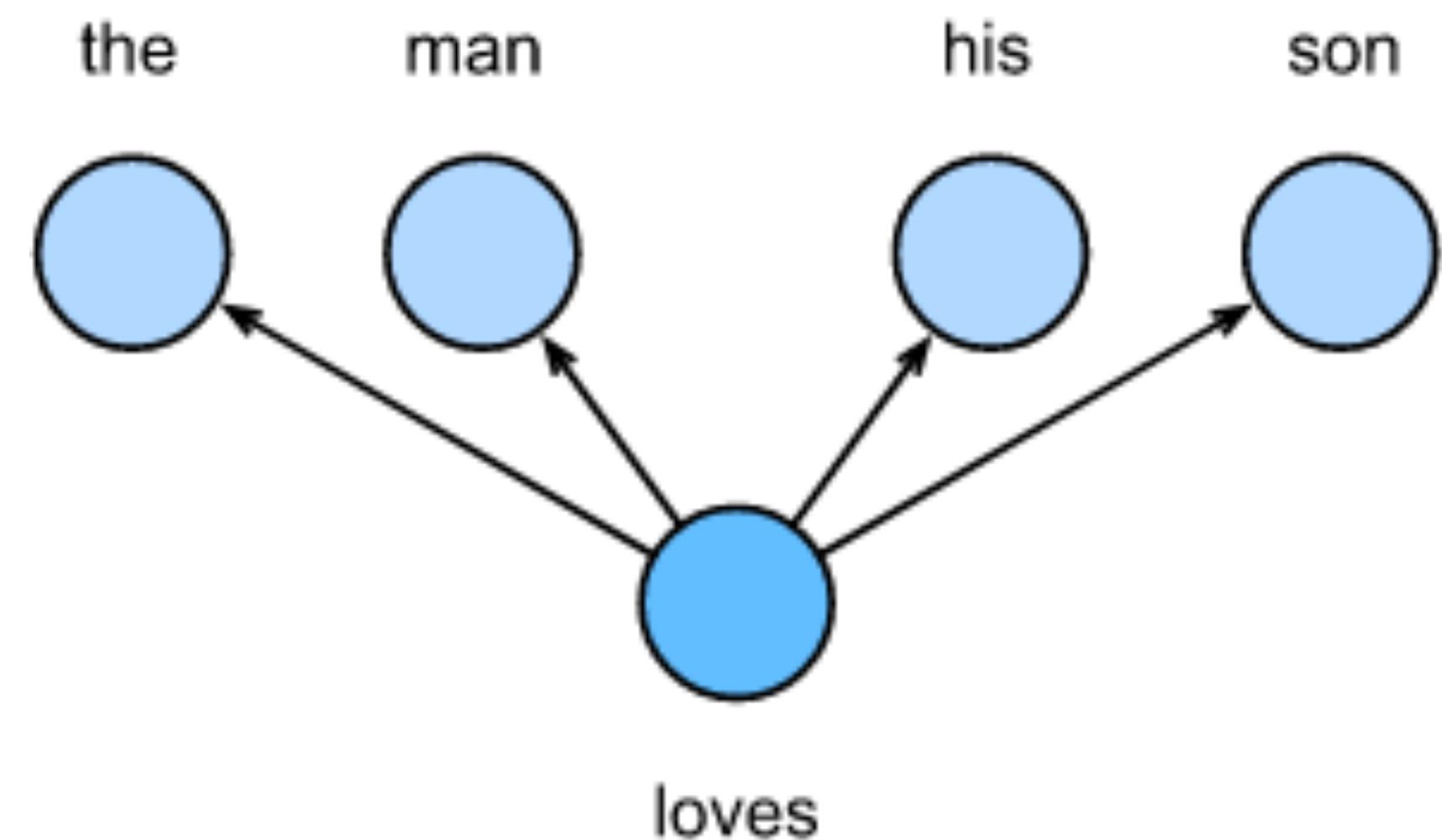
- And $p(w_o \mid w_c) = \frac{\exp(u_o^T v_c)}{\sum_{i \in V} \exp(u_i^T v_c)}$ **cond. probability**, u and v – vector representations.

u_0 - context,

v_c — central target.

- The **likelihood function** of the skip-gram model:

$$\prod_{i=1}^T \prod_{-m \leq j \leq m} P(w^{(t+j)} \mid w^t)$$



Skipgram model training

- Loss function $-\sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \mathbb{P}(w^{(t+j)} | w^{(t)})$
- If we want to SGD it - we need to compute gradient of conditional probability:
- Through differentiation, we can get the gradient from the formula above.
- Any problems?

$$\begin{aligned}\log \mathbb{P}(w_o | w_c) &= \mathbf{u}_o^\top \mathbf{v}_c - \log \left(\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c) \right) \\ \frac{\partial \log \mathbb{P}(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left(\frac{\exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \mathbb{P}(w_j | w_c) \mathbf{u}_j.\end{aligned}$$

Skipgram model training

- Loss function $-\sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log \mathbb{P}(w^{(t+j)} | w^{(t)})$
- If we want to SGD it - we need to compute gradient of conditional probability:
- Through differentiation, we can get the gradient from the formula above:
- Its computation obtains the conditional probability for **all the words in the dictionary** given the central target word w_c
We then use **the same method** to obtain the gradients **for other word vectors**.

$$\begin{aligned}\log \mathbb{P}(w_o | w_c) &= \mathbf{u}_o^\top \mathbf{v}_c - \log \left(\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c) \right) \\ \frac{\partial \log \mathbb{P}(w_o | w_c)}{\partial \mathbf{v}_c} &= \mathbf{u}_o - \frac{\sum_{j \in \mathcal{V}} \exp(\mathbf{u}_j^\top \mathbf{v}_c) \mathbf{u}_j}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \left(\frac{\exp(\mathbf{u}_j^\top \mathbf{v}_c)}{\sum_{i \in \mathcal{V}} \exp(\mathbf{u}_i^\top \mathbf{v}_c)} \right) \mathbf{u}_j \\ &= \mathbf{u}_o - \sum_{j \in \mathcal{V}} \mathbb{P}(w_j | w_c) \mathbf{u}_j.\end{aligned}$$

Negative sampling:

- Given a context window for the central target word w_c , we will treat it as an event for context word w_o to appear in the context window and compute the probability of this event from

$$\mathbb{P}(D = 1 \mid w_c, w_o) = \sigma(\mathbf{u}_o^\top \mathbf{v}_c),$$

- Now we consider maximizing the joint probability $\prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} \mathbb{P}(D = 1 \mid w^{(t)}, w^{(t+j)})$.
- However, the events included in the model only consider positive examples. We need to sample additional negative events (never occurred in the same context) and then:

$$\mathbb{P}(w^{(t+j)} \mid w^{(t)}) = \mathbb{P}(D = 1 \mid w^{(t)}, w^{(t+j)}) \prod_{k=1, w_k \sim \mathbb{P}(w)}^K \mathbb{P}(D = 0 \mid w^{(t)}, w_k).$$

Negative sampling:

- Now we consider maximizing the joint probability

$$\prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} \mathbb{P}(D = 1 \mid w^{(t)}, w^{(t+j)}) .$$

- However, the events included in the model only consider positive examples. We need to sample additional negative K events (never occurred in the same context) and then:

$$\mathbb{P}(w^{(t+j)} \mid w^{(t)}) = \mathbb{P}(D = 1 \mid w^{(t)}, w^{(t+j)}) \prod_{k=1, w_k \sim \mathbb{P}(w)}^K \mathbb{P}(D = 0 \mid w^{(t)}, w_k) .$$

- The logarithmic loss for the conditional probability above is

$$-\log \mathbb{P}(w^{(t+j)} \mid w^{(t)}) = -\log \mathbb{P}(D = 1 \mid w^{(t)}, w^{(t+j)}) - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \mathbb{P}(D = 0 \mid w^{(t)}, w_k)$$

- Here, the gradient computation in each step of the training is no longer related to the dictionary size, but linearly related to K

$$\begin{aligned} &= -\log \sigma \left(\mathbf{u}_{i_{t+j}}^\top \mathbf{v}_{i_t} \right) - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \left(1 - \sigma \left(\mathbf{u}_{h_k}^\top \mathbf{v}_{i_t} \right) \right) \\ &= -\log \sigma \left(\mathbf{u}_{i_{t+j}}^\top \mathbf{v}_{i_t} \right) - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \sigma \left(-\mathbf{u}_{h_k}^\top \mathbf{v}_{i_t} \right) . \end{aligned}$$

Negative sampling:

- The logarithmic loss for the conditional probability above is

$$\begin{aligned} -\log \mathbb{P}(w^{(t+j)} \mid w^{(t)}) &= -\log \mathbb{P}(D = 1 \mid w^{(t)}, w^{(t+j)}) - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \mathbb{P}(D = 0 \mid w^{(t)}, w_k) \\ &= -\log \sigma(\mathbf{u}_{i_{t+j}}^\top \mathbf{v}_{i_t}) - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \left(1 - \sigma(\mathbf{u}_{h_k}^\top \mathbf{v}_{i_t})\right) \\ &= -\log \sigma(\mathbf{u}_{i_{t+j}}^\top \mathbf{v}_{i_t}) - \sum_{k=1, w_k \sim \mathbb{P}(w)}^K \log \sigma(-\mathbf{u}_{h_k}^\top \mathbf{v}_{i_t}). \end{aligned}$$

- Here, the gradient computation in each step of the training is no longer related to the dictionary size, but linearly related to K

- **Key idea:** sample additional negatives and learn your positive probabilities «opposite to» them.

So what? (Synonyms)

```
get_similar_tokens('chip', 3, glove_6b50d)
```

```
cosine sim=0.856: chips  
cosine sim=0.749: intel  
cosine sim=0.749: electronics
```

```
get_similar_tokens('baby', 3, glove_6b50d)
```

```
cosine sim=0.839: babies  
cosine sim=0.800: boy  
cosine sim=0.792: girl
```

```
get_similar_tokens('beautiful', 3, glove_6b50d)
```

```
cosine sim=0.921: lovely  
cosine sim=0.893: gorgeous  
cosine sim=0.830: wonderful
```

- *get_similar_tokens* – top-K words by cosine measure to the target word;
- *glove_6b50d* – glove model on some common corpora (Wikipedia?) with 6B of words and vector dimension equals to 50;

So what? (2) (Finding Analogies)

```
get_analogy('man', 'woman', 'son', glove_6b50d)
```

```
'daughter'
```

“Capital-country” analogy

```
get_analogy('bad', 'worst', 'big', glove_6b50d)
```

```
'biggest'
```

“Adjective-superlative adjective”
analogy

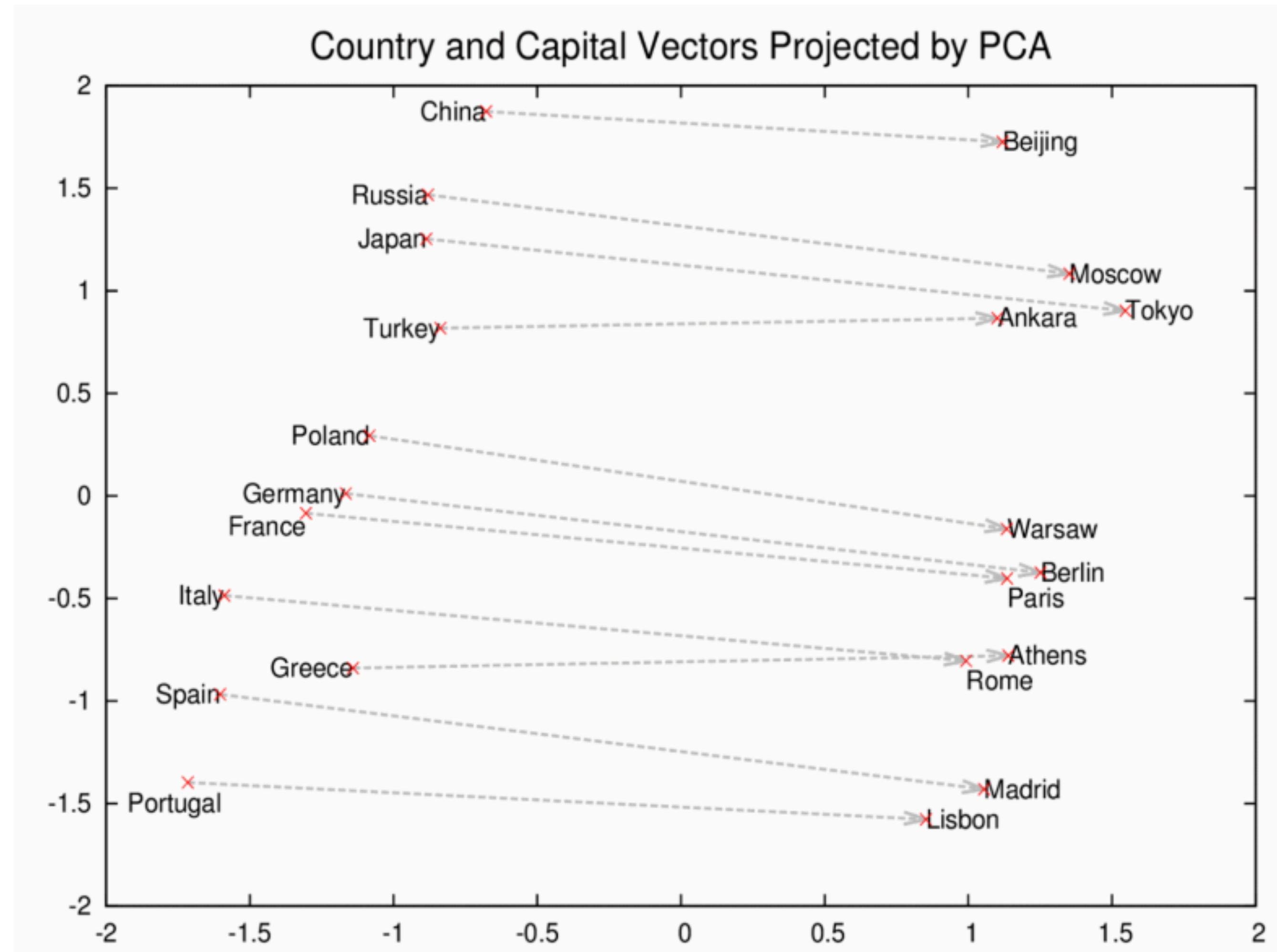
```
get_analogy('do', 'did', 'go', glove_6b50d)
```

```
'went'
```

“Present tense verb-past
tense verb” analogy

- And it's only $x = \text{vec}(c) + \text{vec}(b) - \text{vec}(a)$
- And then top word for x

So what? (country-capital)



Based on Wikipedia training corpora

Any problems?

- Out-of-vocabulary
- How we can train it? How big our doc's collection should be?
- Stop, firstly we talk about **text** and word2vec is about **words**

Any problems?

- Out-of-vocabulary

Yeah, it's true. But there are few extensions; (fastText)

- Stop, firstly we talk about **text** and word2vec is about **words**

Ok, average it; Or average with weights; Or do not average and learn some averaging embedding; (look to BERT model)

- How we can train it? How big our doc's collection should be?

Really big; Starting from 10+M of symbols; Use pertained vectors;

Word2Vec in RecSys

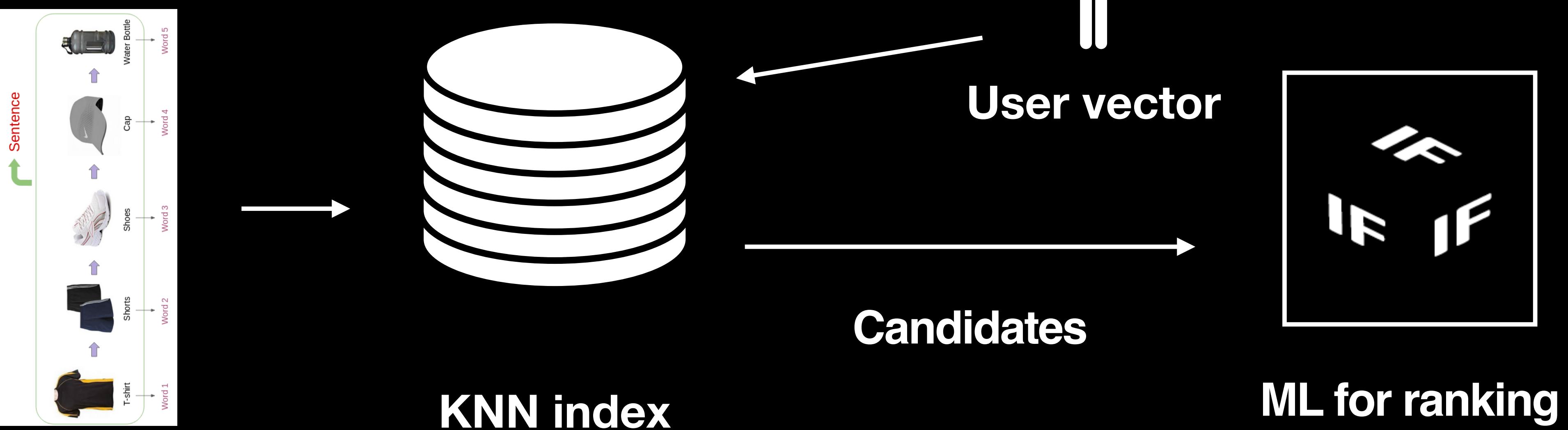
- Negative sampling
- Next item in sequence as a task
- Good initialisation for «Embedding layer»
- Filter your datasets, please



Word2Vec in RecSys

- Negative sampling
- Next item in sequence as a task
- Good initialisation for «Embedding layer»
- Filter your datasets, please

Problem?



Word2vec offline model with embeddings

W2V:

- **No user** vector:
 - Sum all items user interacted with
 - Sum with weights ~ activity
 - Sum with TF-IDF
 - Cluster and select medoid...
 - Feed to an attent...
- **Pros:**
 - Easy to realise
 - No need to limit users interaction
 - Good at similars
 - Session-based
 - Consume CPU , not GPU
- **Cons:**
 - Nothing except interactions
 - Cold start with items
 - No user vector

W2V

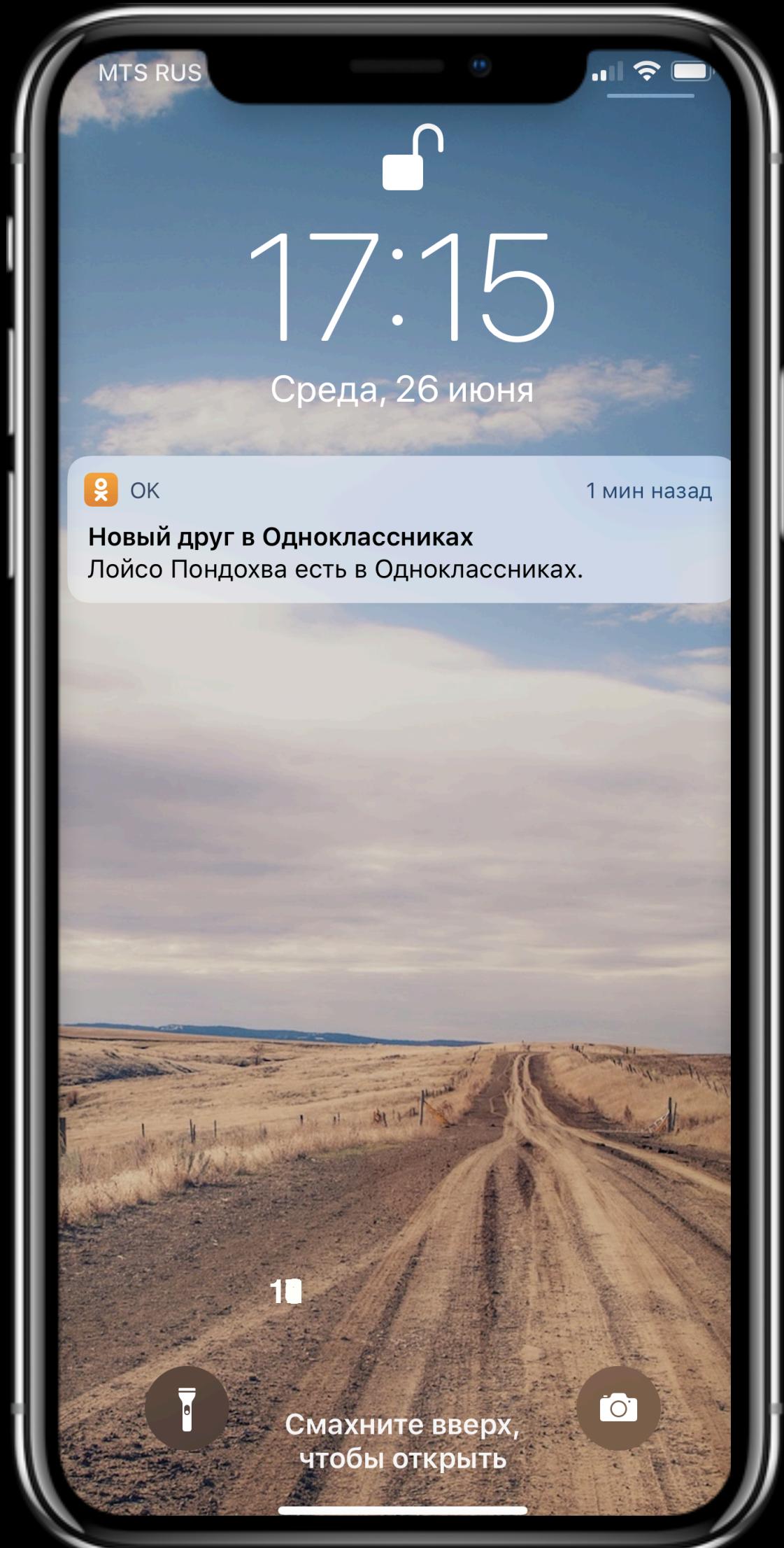
Discussion notes

- Items are limited by number of interactions – how to overcome it?
- Intrinsic evaluation. What to do?



- Rare/popular items lost their similars quality, why?

РУМК:

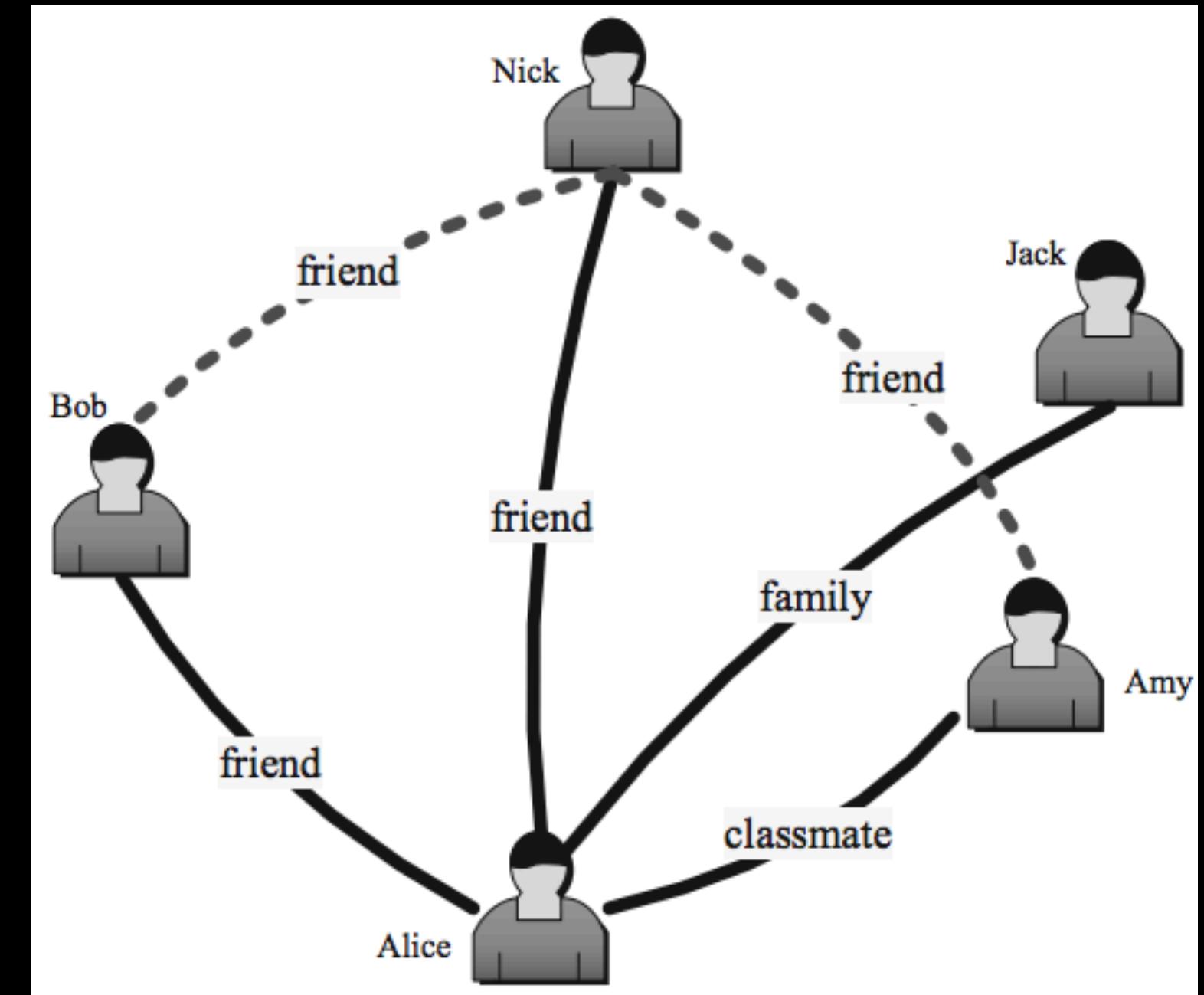


The main screen of the 'OneClass' app displays several sections of user profiles:

- Рекомендую** (Recommended): Shows a profile picture of a white spider with a flower crown. Below it, the name 'Наида Тагирова' and '4 общих друга' (4 mutual friends) are listed, along with a 'Добавить в друзья' (Add to friends) button.
- Илья Родионов**: Shows a profile picture of a man in sunglasses. Below it, the name 'Илья Родионов' and '11 общих друзей' (11 mutual friends) are listed, along with a 'Добавить в друзья' (Add to friends) button.
- Вы можете знать этих людей** (You may know these people): A grid of eight profile cards, each with a 'Дружить' (Friendship) button below it:
 - Максим Овечкин**: Profile pic, 39 mutual friends.
 - Vladislav Dolganov**: Profile pic, 28 mutual friends.
 - Вадим Баранов**: Profile pic, 41 mutual friend.
 - Тимур Насрединов**: Profile pic, 67 mutual friends.
 - Vadim Tsesko**: Profile pic, 63 mutual friends.
 - Никита Макаров**: Profile pic, 68 mutual friends.
 - Sergey Alekseev**: Profile pic, 46 mutual friends.
 - Михаил Марюфич**: Profile pic, 26 mutual friends.

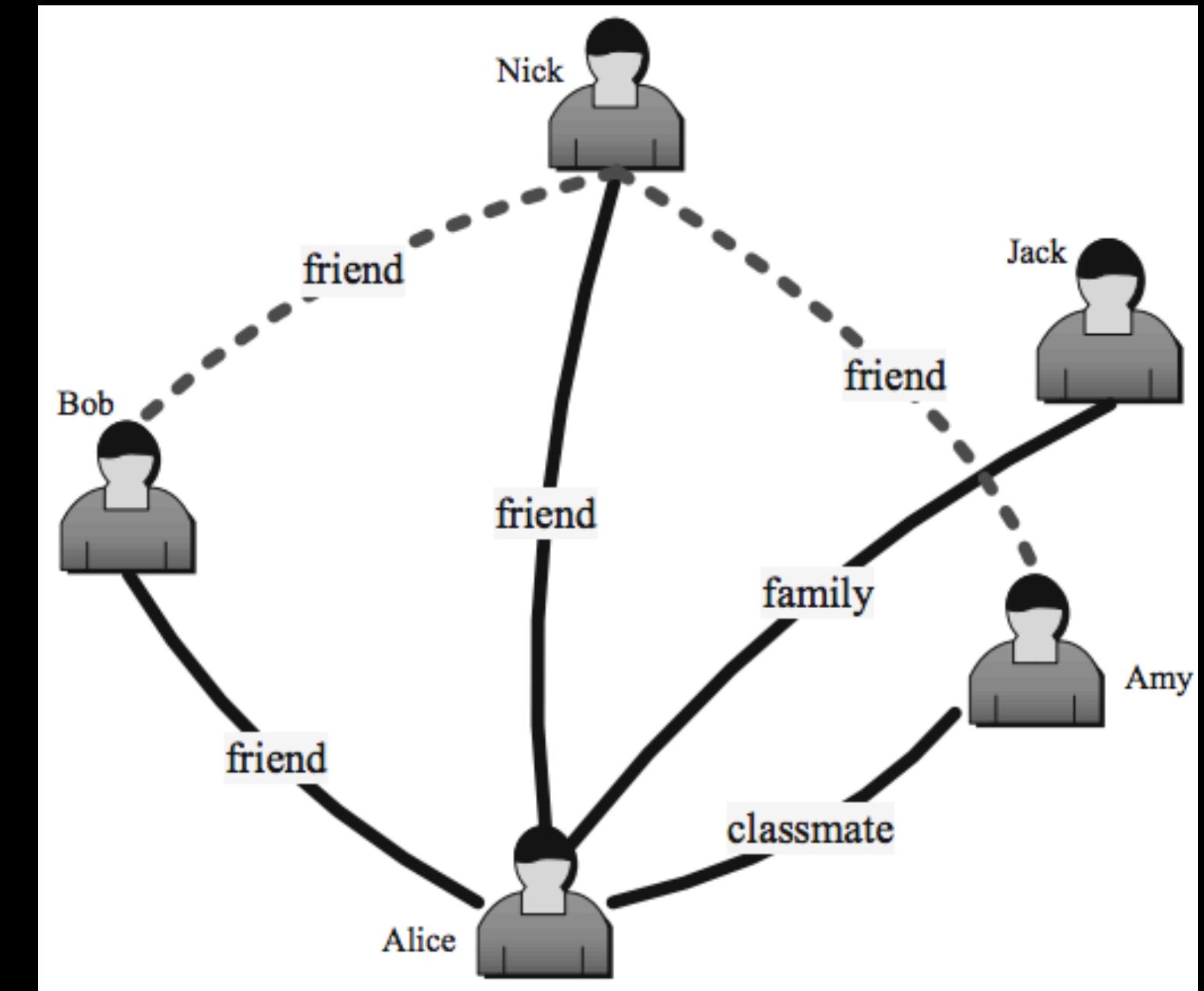
Link prediction(recovery):

- Treat link prediction as binary classification.
- Each unconnected pair – corresponds to an instance with features and label.
- If there is a potential link – instance labeled as «positive», otherwise – negative.



Link prediction(recovery):

- Treat link prediction as binary classification.
- Each unconnected pair – corresponds to an instance with features and label.
- If there is a potential link – instance labeled as «positive», otherwise – negative.

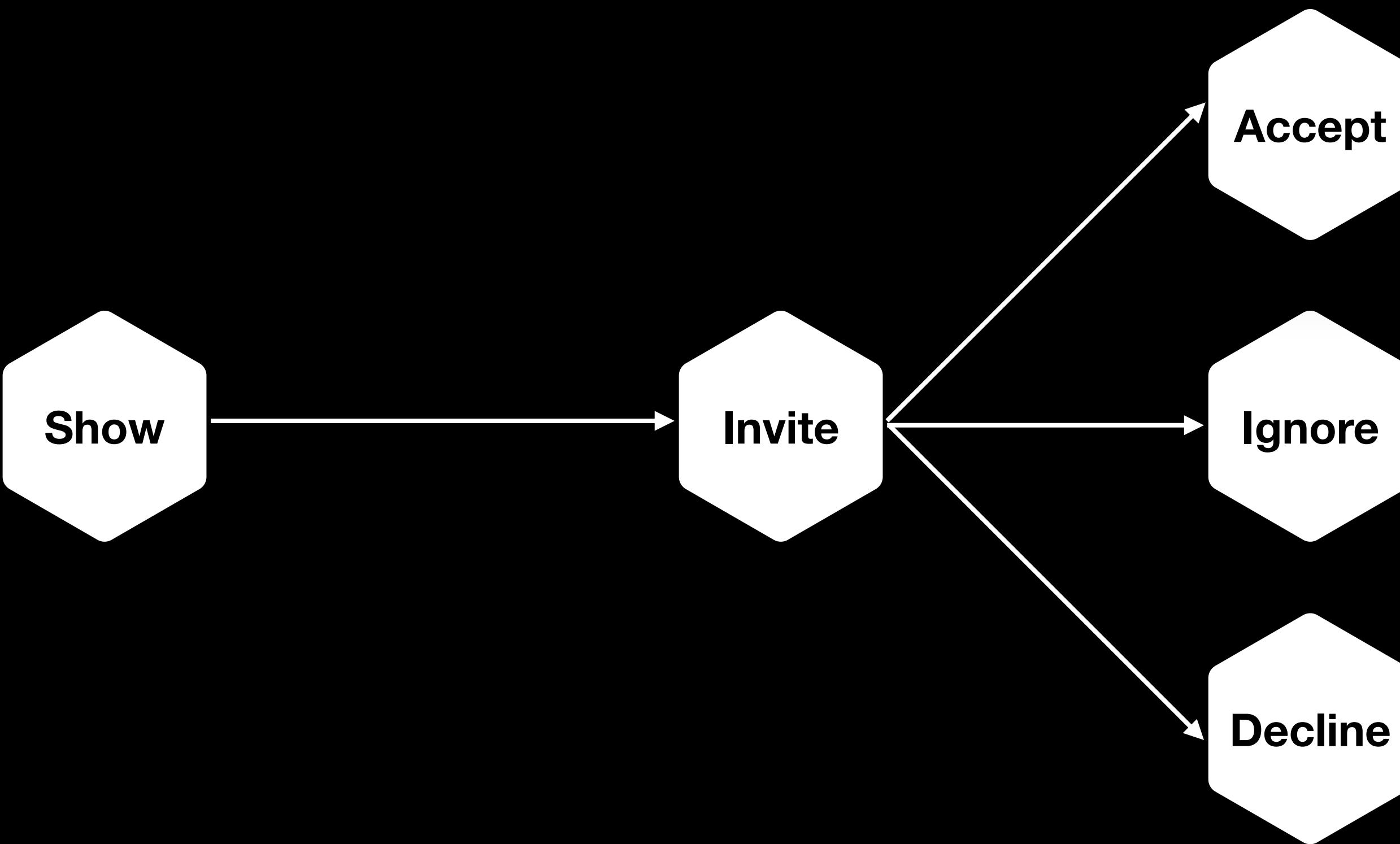


PROBLEMS?

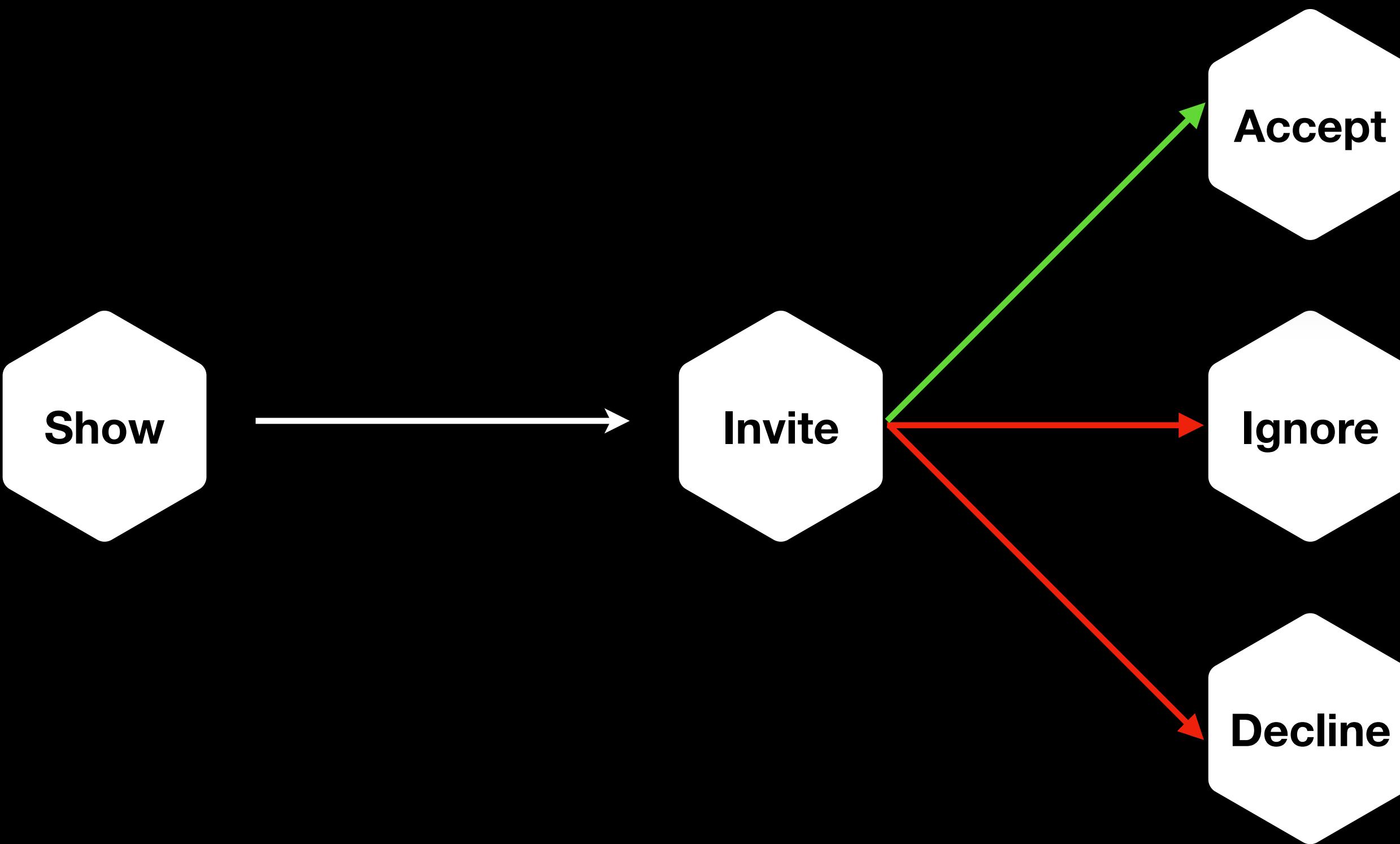
Link prediction(recovery):

- Treat link prediction as **binary** classification.
- **Each unconnected pair** – corresponds to an instance with features and label.
- If **there is a potential link** – instance labeled as «positive», otherwise – negative.
- What should we use as a **target**?
- Should we score **each unconnected pair** of nodes in 200M+ nodes graph?
- So what should **we exactly do** when «there is a potential link»?

Target



Target

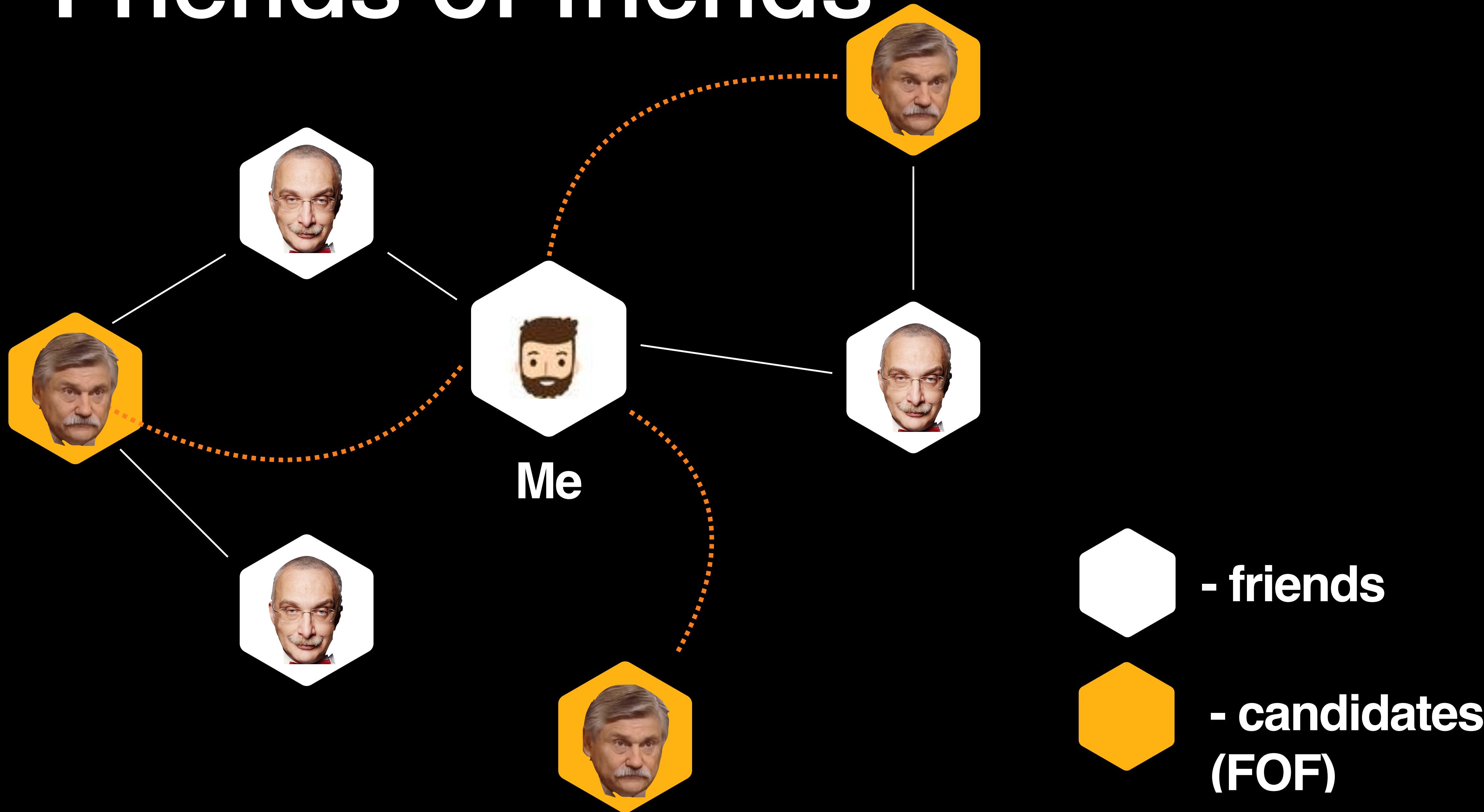


Recommendation System schema:



- **Matrix factoriation + K-NN**
- **Recently acted**
- **Friends-of-friends**
- **Reverse index**
- ...
- **Embeddings dot**
- **Neural Networks**
- ★ **Machine learning model
(xgboost ^^)**

Friends of friends

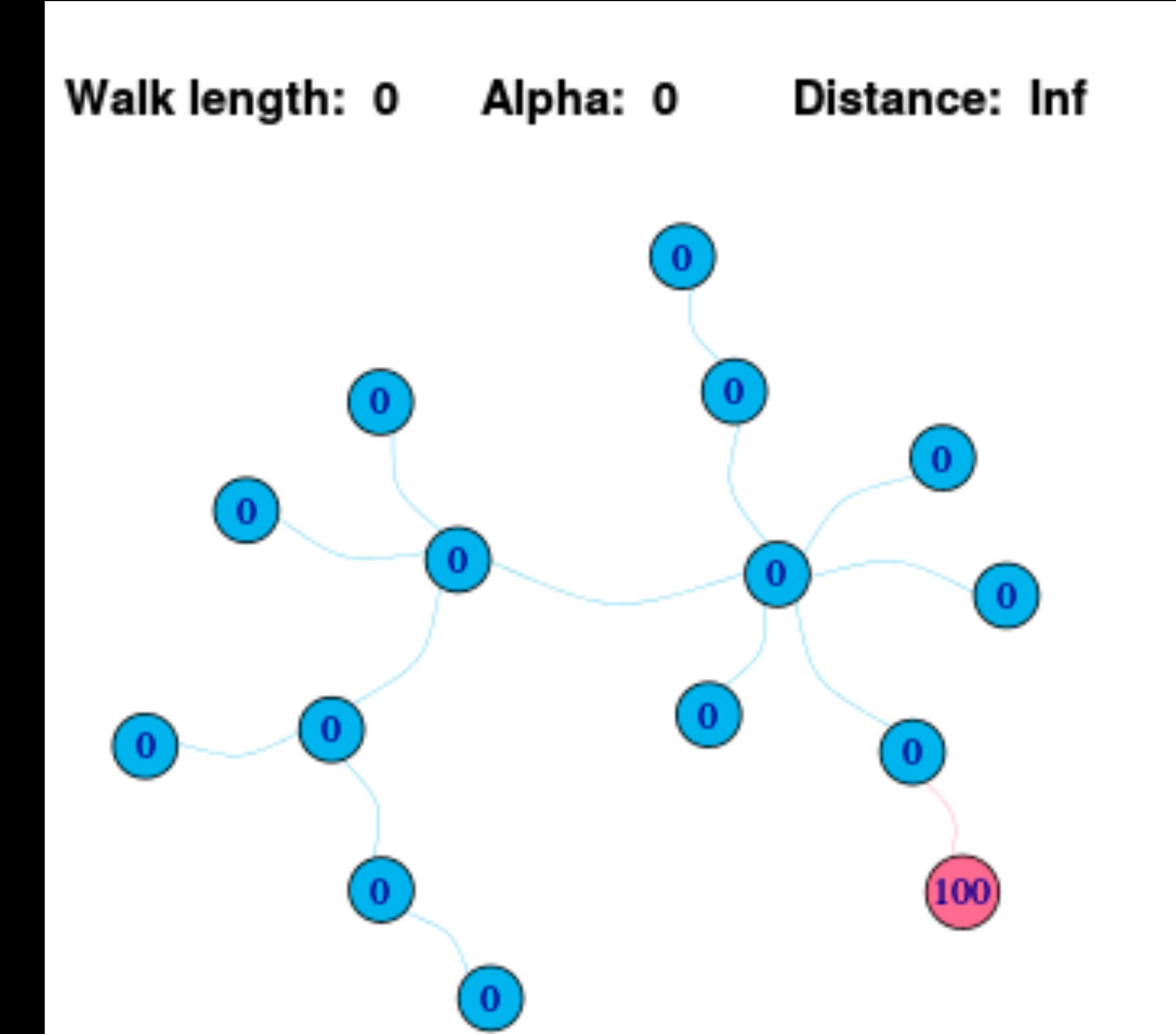


What's more?

- Last active interaction
- Phonebook
- Reverse phonebook
- Together on photo
- Graph embeddings... ?

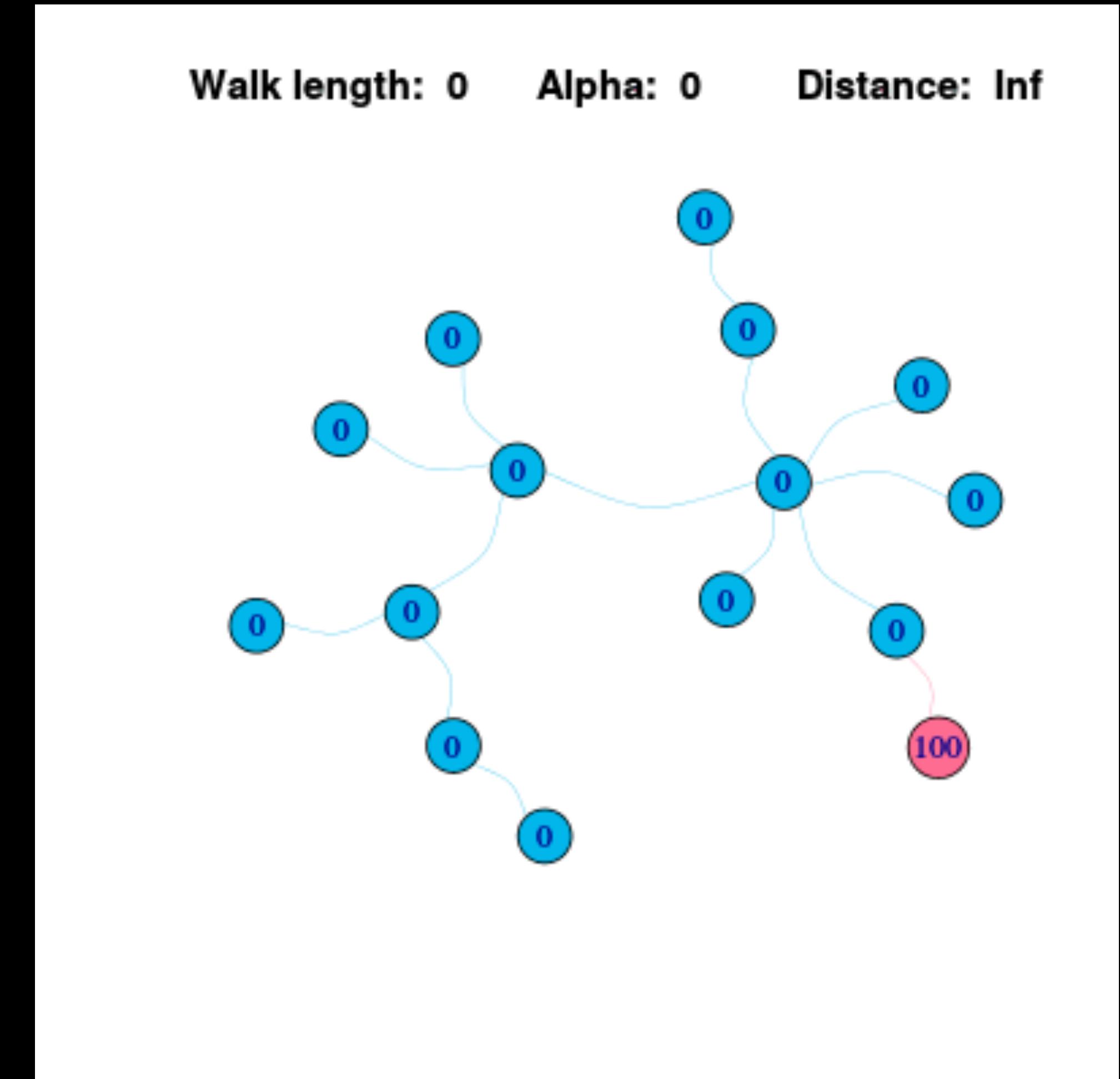
PageRank: random walks

- Where does it spread first?
- How far does it spread?
- Will the vertices, close to the red vertex obtain more information then those far away?
- Will the information continue to go back and forth forever or will we reach a stable distribution eventually?



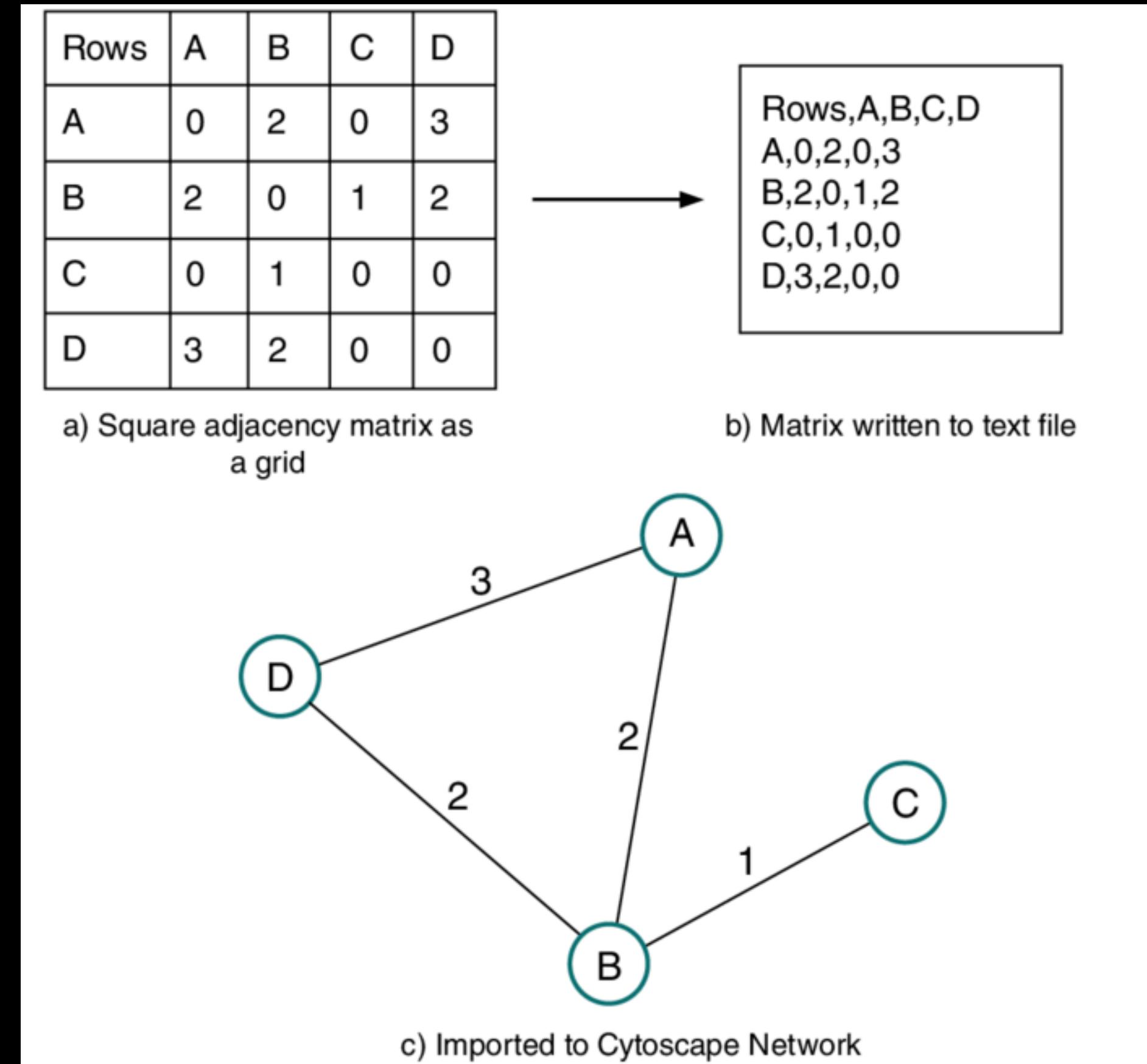
PageRank: random walks

- Goes from red through random edge to next adjacent vertex
- If we repeat thus experiment infinite number of times - will it converge?
- Let assume A – matrix of transition probabilities, x – initial distribution.
- 1st step $x' = Ax$
- 2nd step $x'' = Ax'$
- n 'th step: $x^n = A x^{(n-1)}$



PageRank: random walks

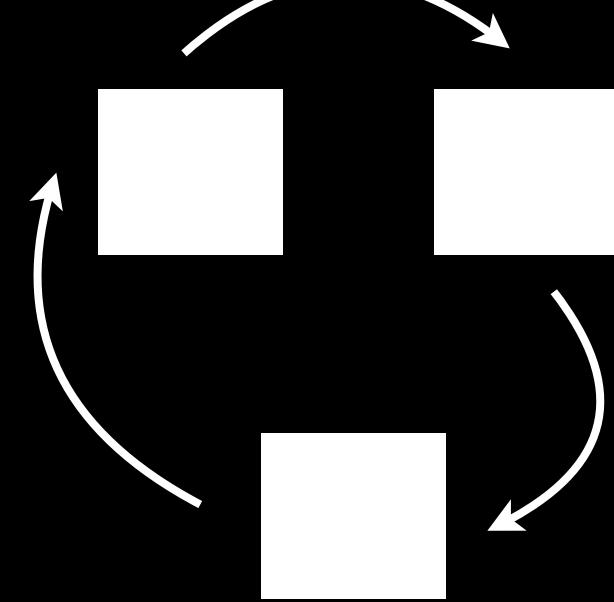
- 1st step $x' = Ax$
2nd step $x'' = Ax'$
 n 'th step: $x^n = A x^{(n-1)}$
- Convergence: $x^n = A * x^{(n-1)} \dots$
- $Ax = \lambda x$ doesn't look similar?
- Final distribution of weights - is Eigenvector of matrix A with eigenvalue lambda = 1.
- Do you feel any **problems**?
- But wait ... how do we know that A has an eigenvalue of 1, so that such a vector p exists? And even if it does exist, will it be unique (well-defined)?



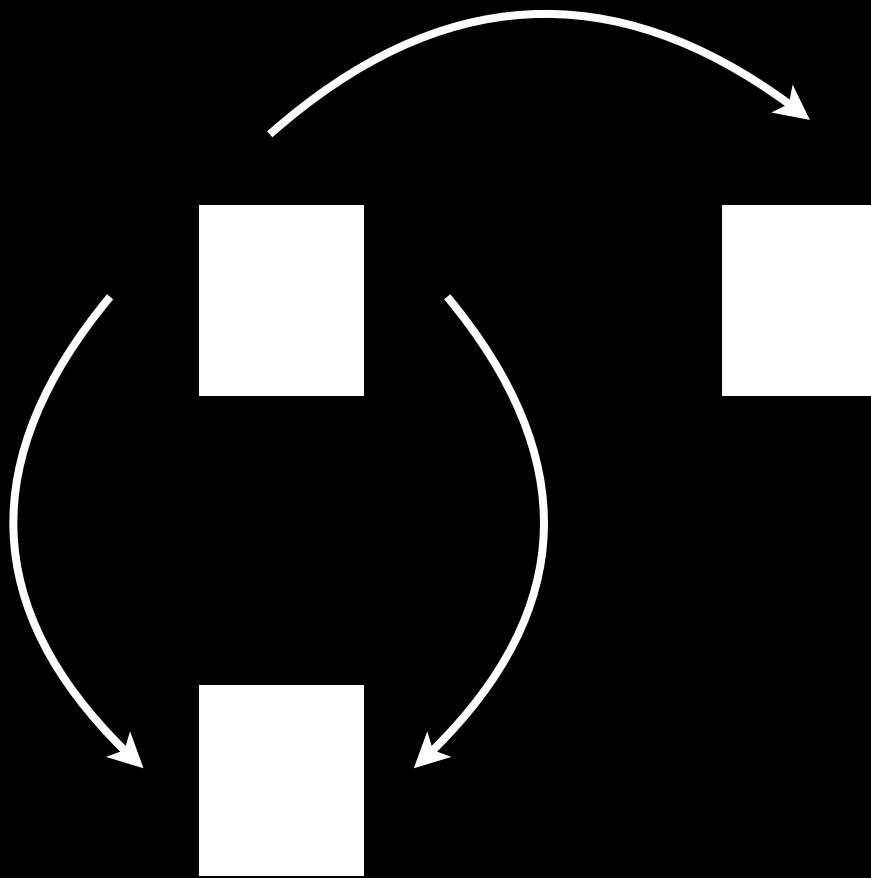
PageRank: as Markov chain

- Markov chain - random process over states $s_1, s_2 \dots s_n$ defined by transition matrix P $n \times n$.
- Let denote $p^{(0)}$ as initial distribution vector. $p^{(n)} = P * p^{(n-1)}$
- A **stationarity distribution** is probability vector $p = P * p$ (oops)
- If the Markov chain is **strongly connected**, meaning that any state can be reached from any other state, then the stationary distribution **p exists and is unique**. Furthermore, we can think of the stationary distribution as the proportions of visits the chain pays to each state after a very long time (the ergodic theorem):
- P - problems?

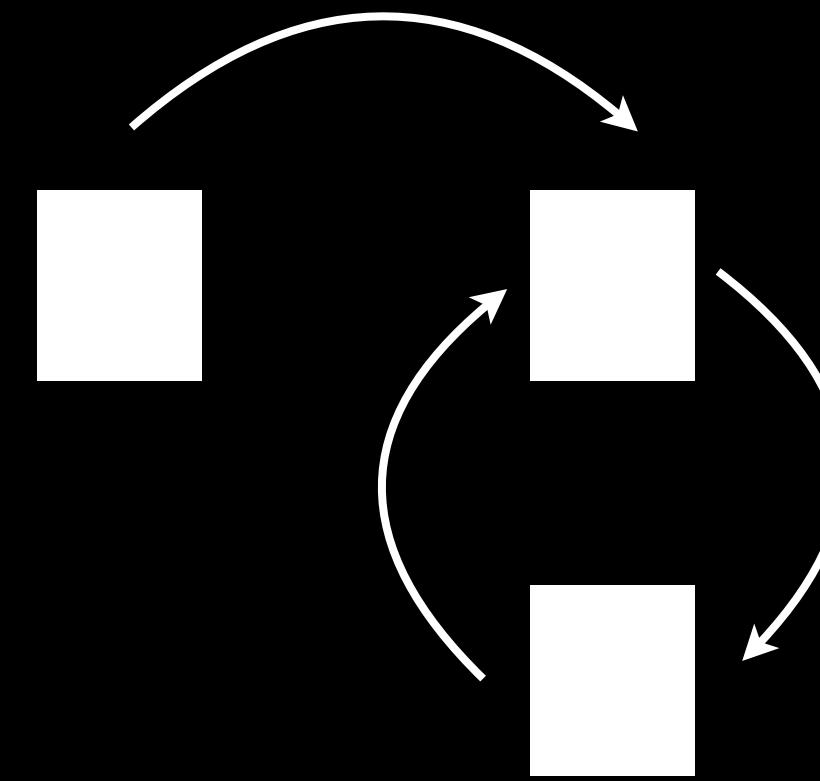
PageRank as Markov chain:



Disconnected components



Dangling links

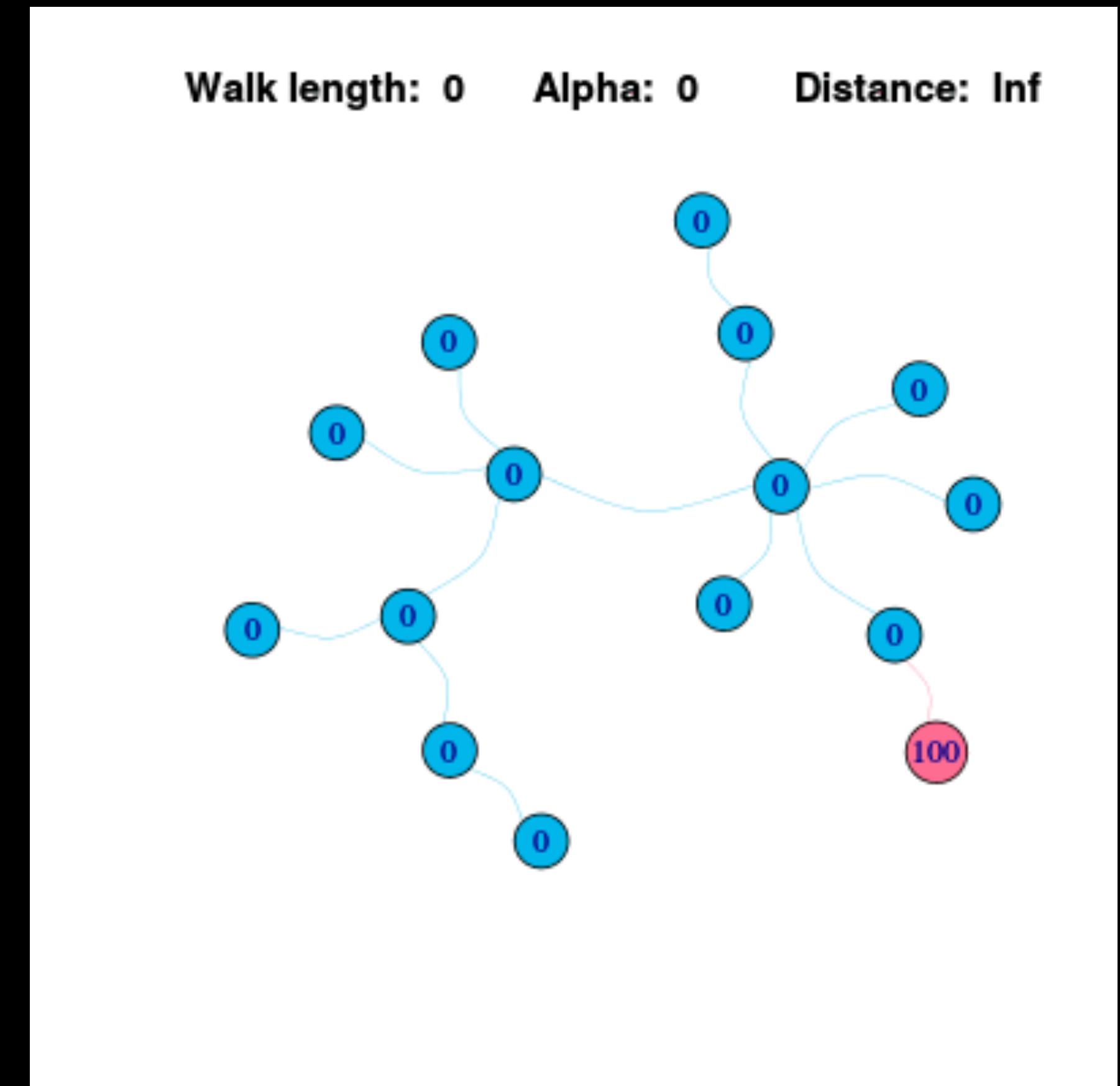


Loops

- Actually, even for Markov chains that are not strongly connected, a stationary distribution always exists, but may **nonunique**.

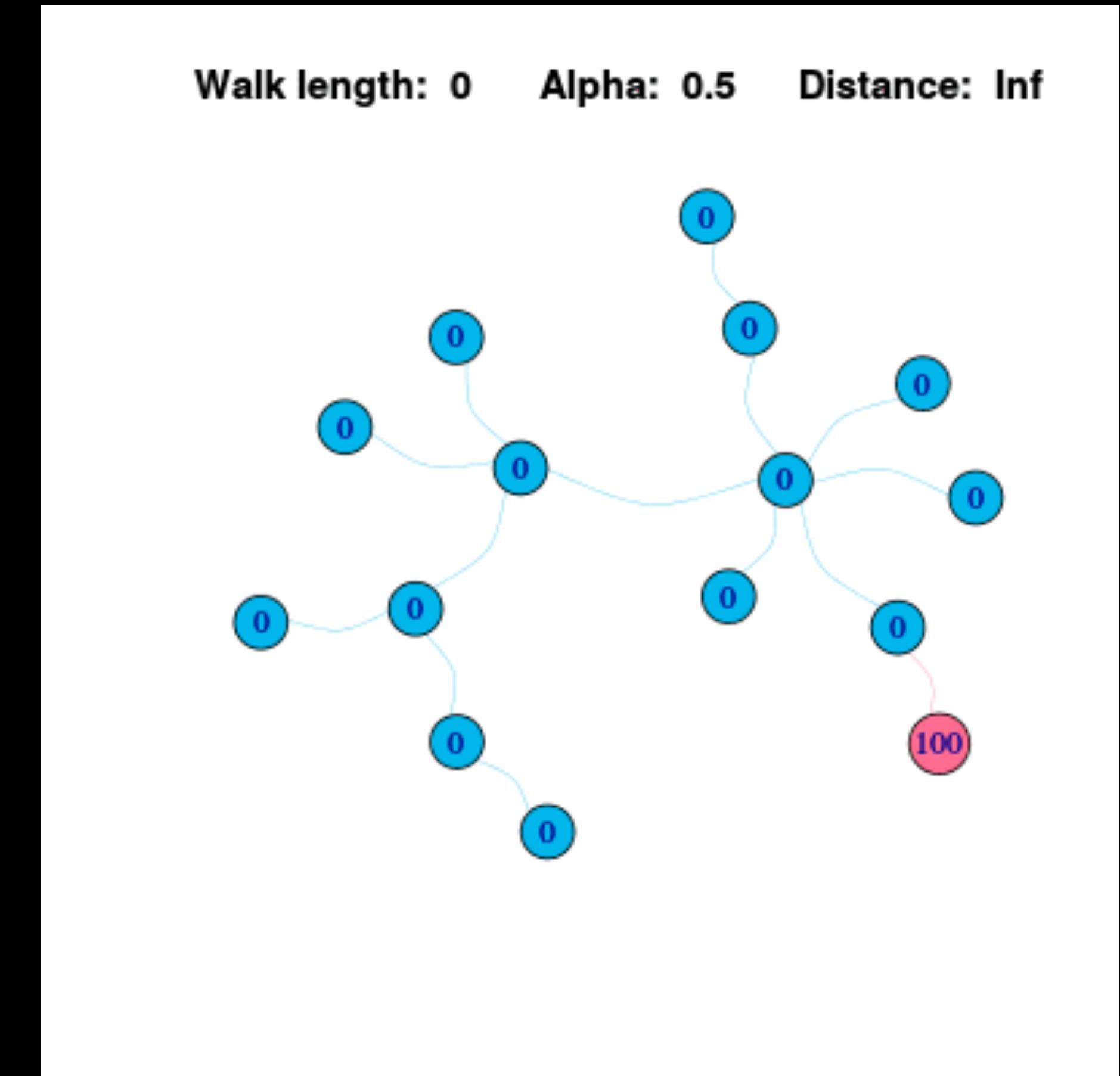
Lazy random walks

- The idea of **lazy** random walks: we allow the random walkers to remain on a vertex with probability 1/2.
- Leads to $x' = 1/2(A + I) * x$. There are some more complicated formulas for eigenvectors.
- Edges with an adjacent vertex that has a high degree tend to be colored red. These are the vertices which contain a significantly larger part of information than the rest of the network.
- No difference where information **originates from**. This means that high-degree vertices will contain proportionally more information than low-degree vertices.



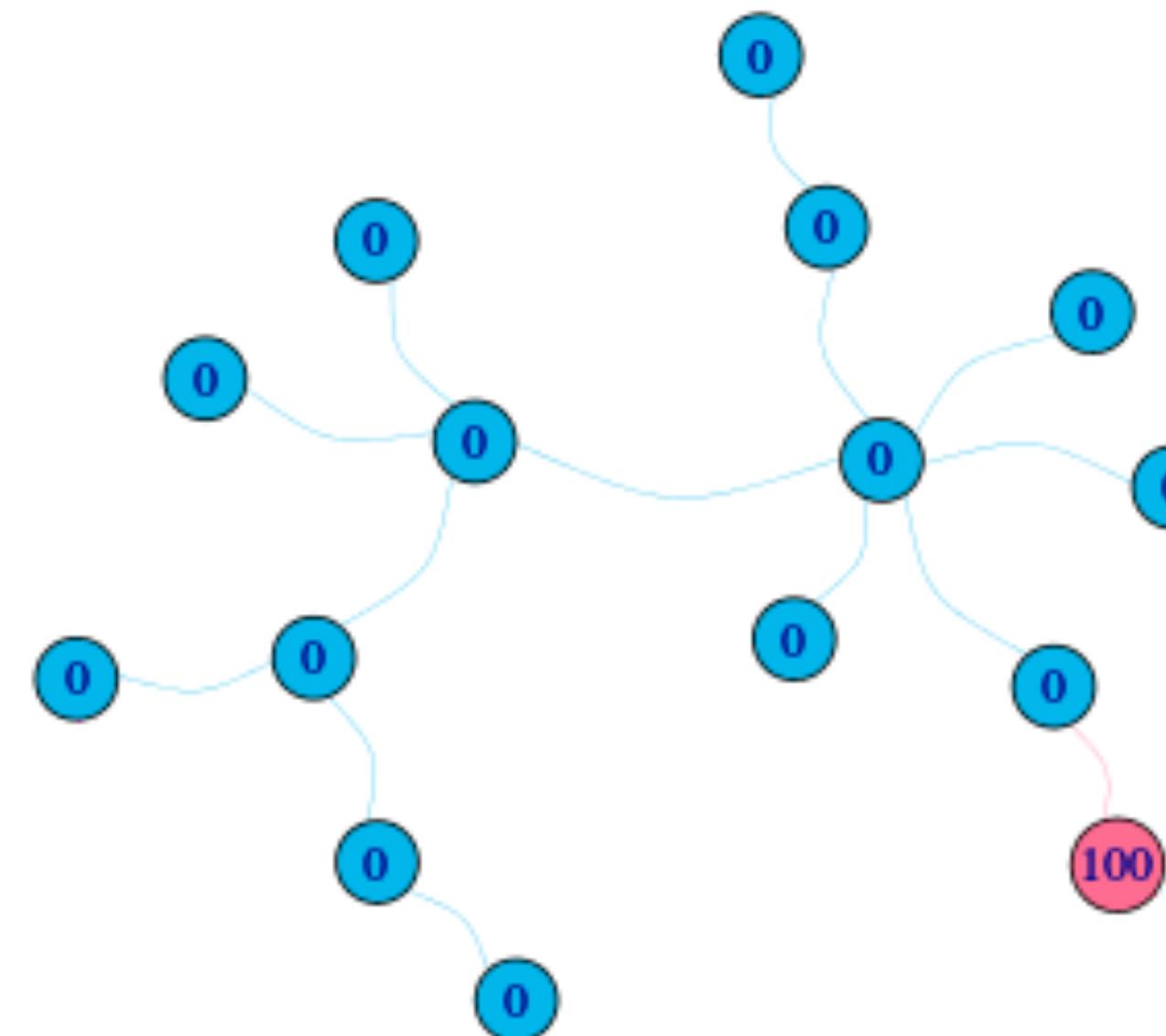
Personalized PageRank:

- Let introduce α as probability to teleport back to source vertex. It tends to make our distribution more «**centered**» on target vertex
- $x' = (1 - \alpha) * Ax + \alpha * E$, E - with «target» distribution, «one-hot» with 1 at target vertex, for example.
- This schema also known as «personalized page rank». We may interpret it as «walks» from home with **avg length: 1/alpha**



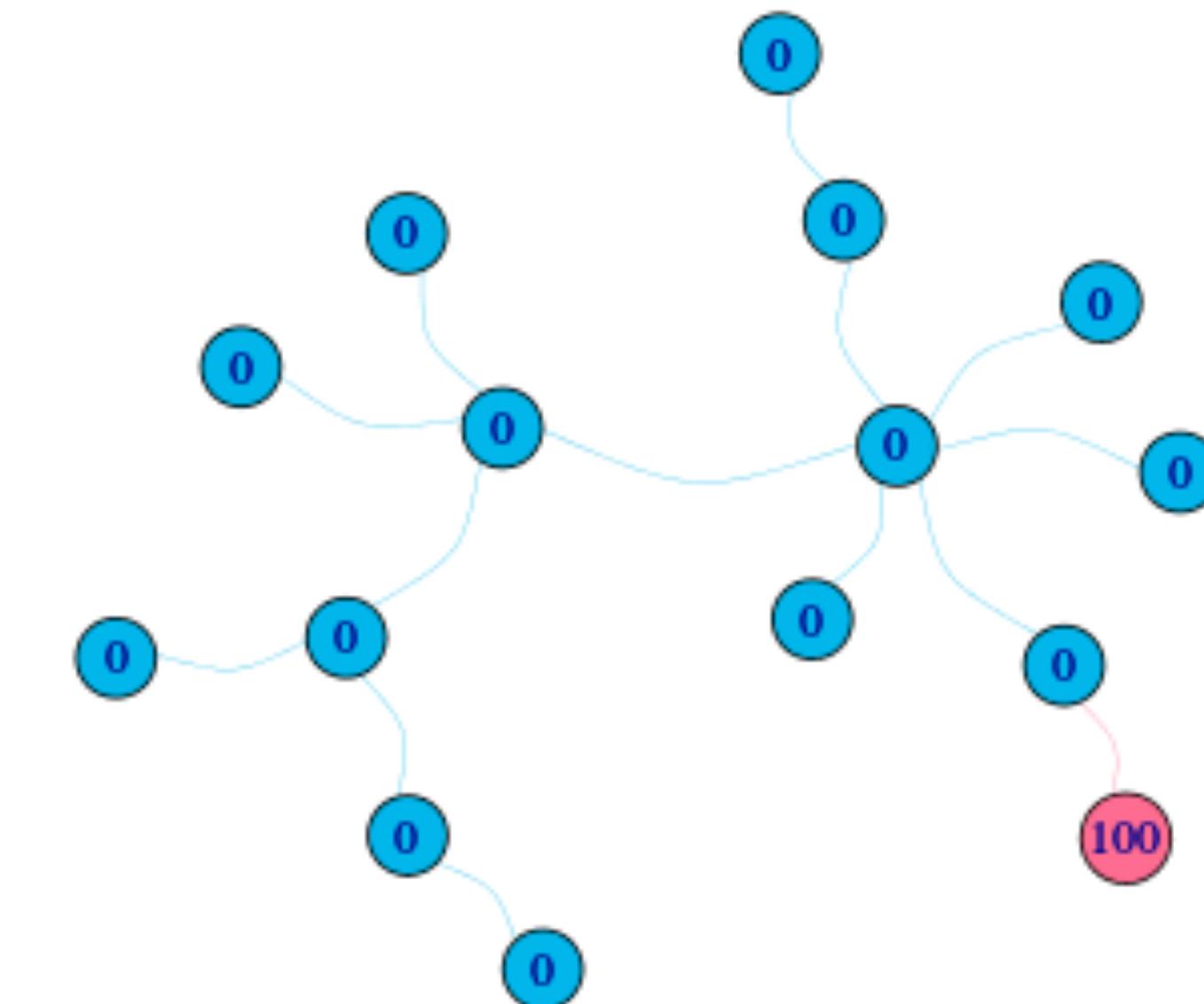
Personalized PageRank:

Walk length: 0 Alpha: 0.1 Distance: Inf



Alpha = 0.1

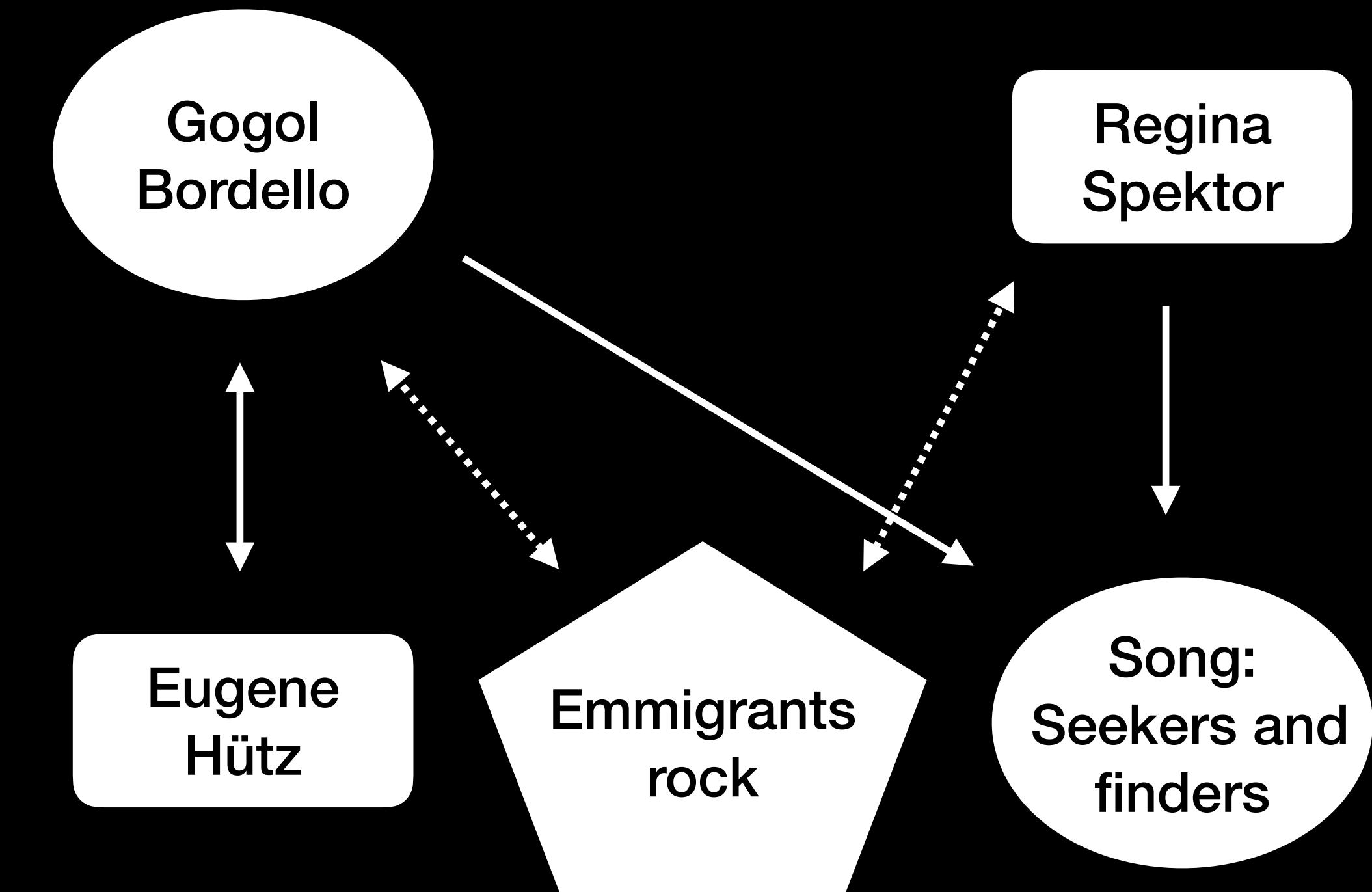
Walk length: 0 Alpha: 0.5 Distance: Inf



Alpha = 0.5

PageRank applications

- Google's first score algorithm
- PYMK as candidates generator
- Not only PYMK, but any other graph-representable product. Music recommender, for example:
<https://arxiv.org/pdf/1310.7428.pdf>



PageRank tips and **tricks**

- Do you really wanna search for eigenvalues thus matrixes for every user?
 $x' = (1 - \alpha) * Ax + \alpha * E$, in prod we just launches few raw random walks and collect results.
- There are a lot of hyperparameters. For example, different weights for different type of edges. I know no method to tune it but **grid search and a/b tests**.
- Good starting point to bring outer worlds knowledge to your model:
Look at SPARQL

<https://www.w3.org/TR/sparql11-query/>

Query language for semantic web DB based on wikipedia.

PageRank discussion:

- Random walks instead of eigenvectors in runtime
- A lot of hyperparameters (different types of edge/vertex) and etc.
- **No embeddings** space, only scores and candidates
- ...?

BigGraph by Facebook:

- Facebook graph-scale system to train graph embeddings
- $G = (V, R, E)$ - nodes, set of relations, edges
- $f(\theta_s, \theta_r, \theta_d) = sim(g_s(\theta_s, \theta_r), g_d(\theta_d, \theta_r))$ - score function.

$\theta_s, \theta_r, \theta_d$ – embeddings of source, relation and destination

g_s, g_d – mapping operators to introduce different types of mapping for relations

- Negative sampling + distributed engine

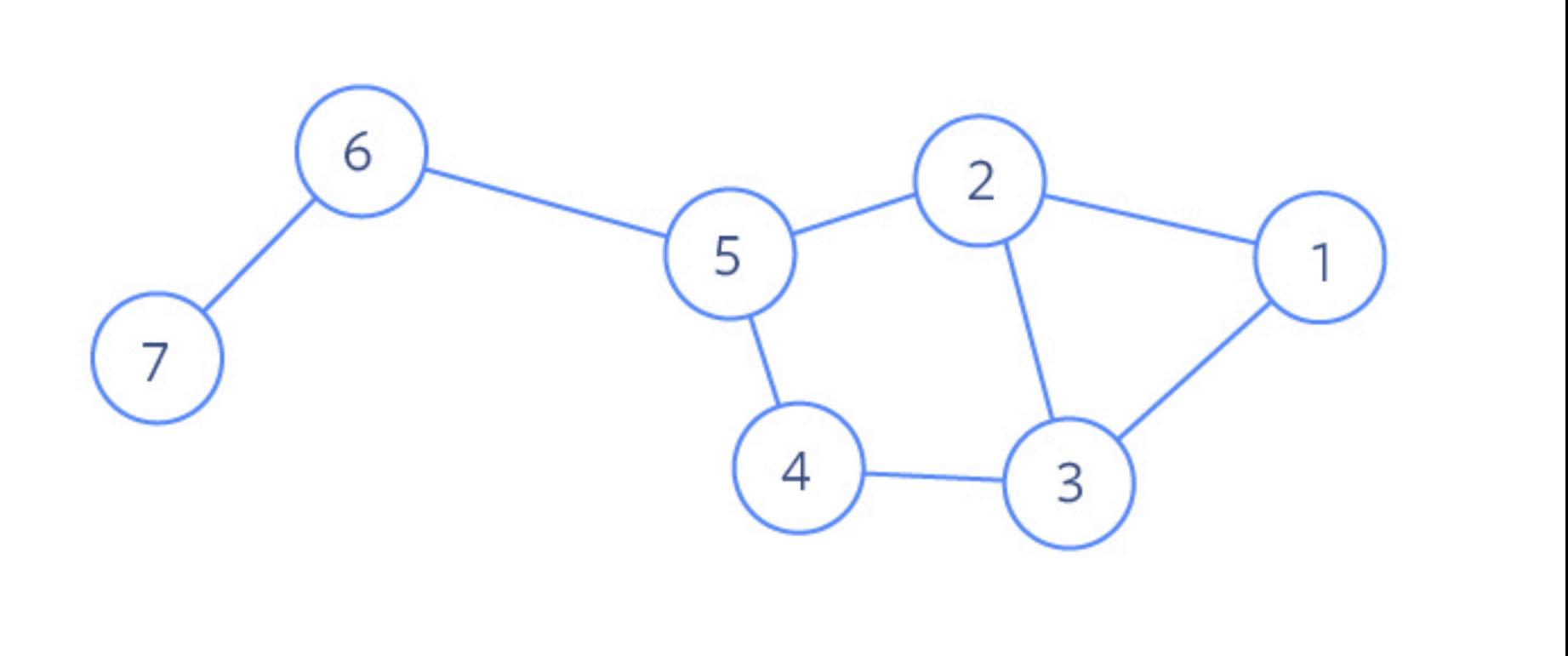
- Imagine:

$$g_s(\theta, \theta_r) = g_d(\theta, \theta_r) = \theta$$

$$sim(\theta_s, \theta_d) = cos(\theta_s, \theta_d)$$

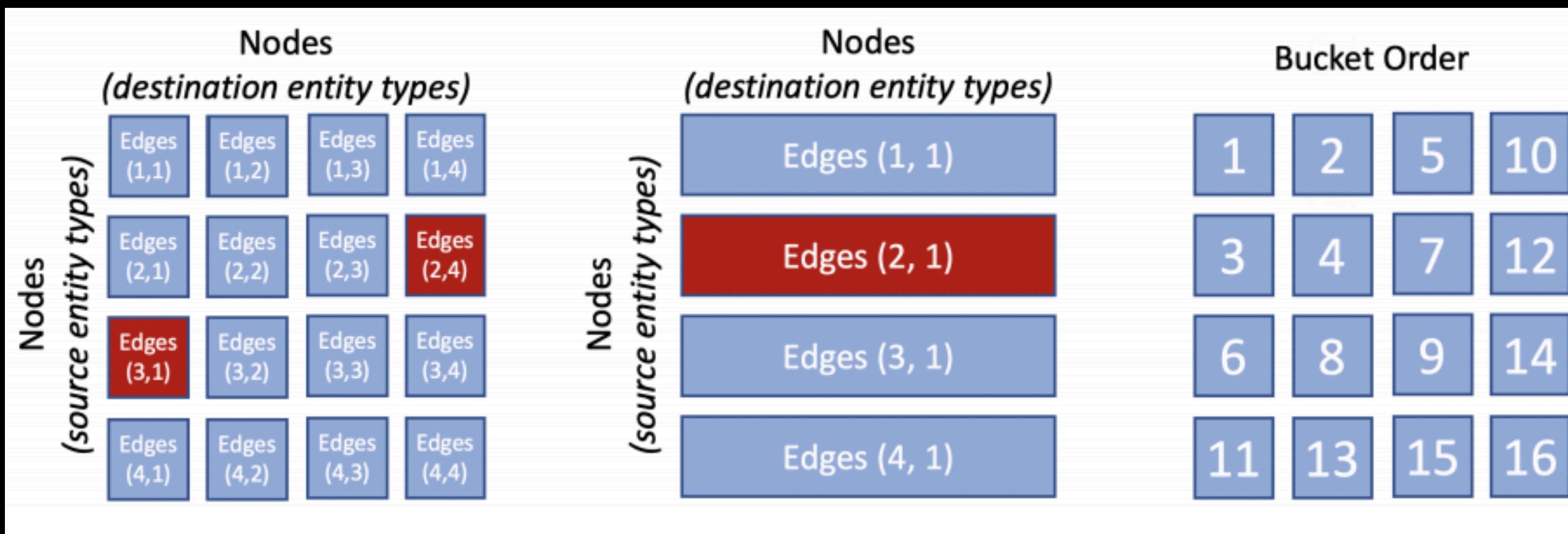
^
^

formally equals word2vec



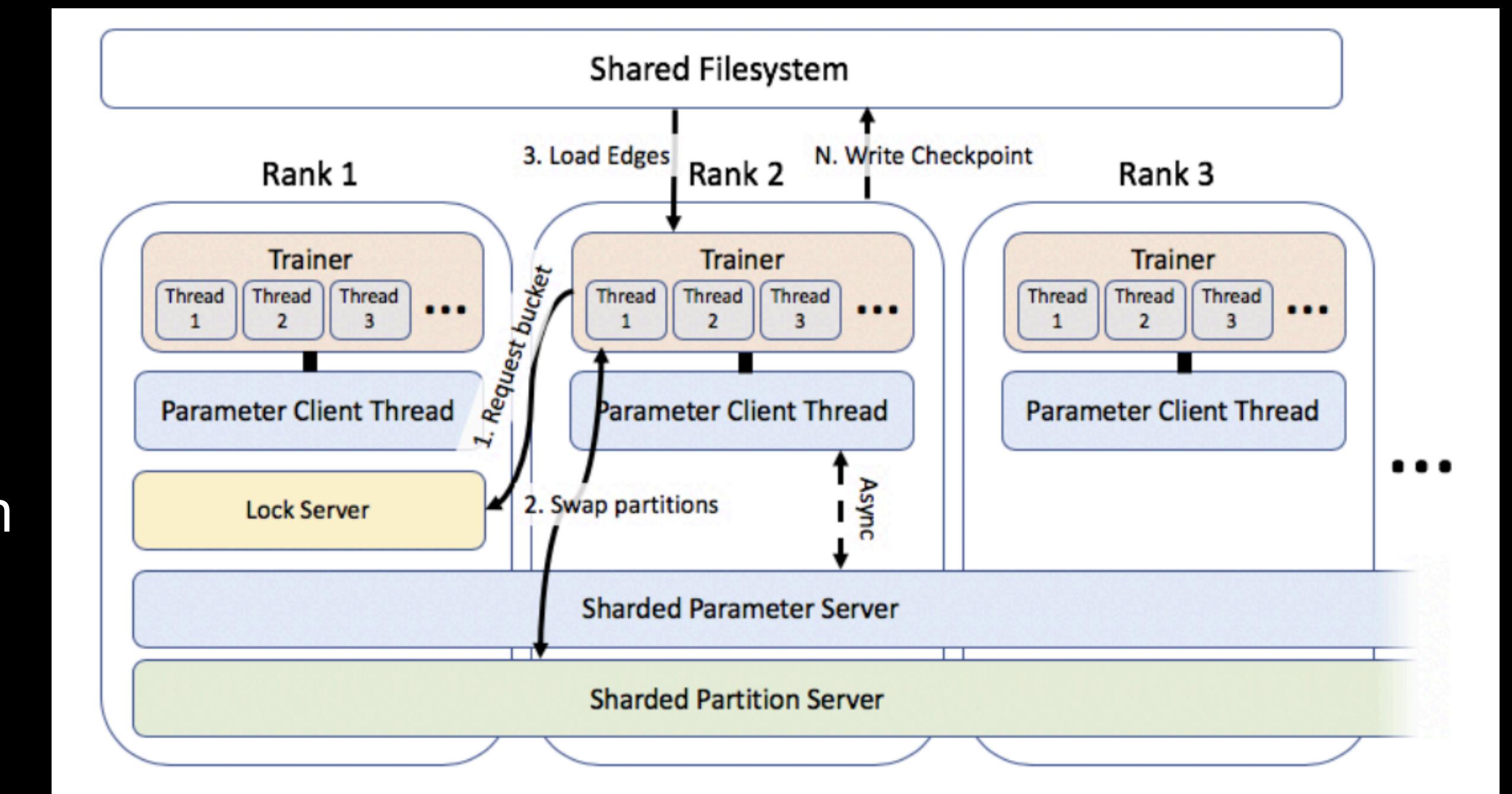
BigGraph:

- **Left:** nodes are divided into P partitions that are sized to fit in memory. Edges are divided into buckets based on the partition of their source and destination nodes. In distributed mode, multiple buckets with non-overlapping partitions can be executed in parallel (red squares).
- **Center:** Entity types with small cardinality do not have to be partitioned; if all entity types used for tail nodes are unpartitioned, then edges can be divided into P buckets based only on source node partitions.
- **Right:** the ‘inside-out’ bucket order guarantees that buckets have at least one previously-trained embedding partition. Empirically, this ordering produces better embeddings than other alternatives (or random



BigGraph:

- Rank 2 Trainer performs for the training of one bucket.
- Trainer requests a bucket from **the lock server** on Rank 1, which locks that bucket's partitions.
- The trainer then saves any partitions that it is no longer using and loads new partitions that it needs to and from the sharded partition servers, at which point it can release its old partitions on the lock server.
- Edges are then loaded from a shared filesystem, and training occurs on multiple threads without inter-thread synchronization(Recht et al., 2011).
- In a separate thread, a small number of shared parameters are continuously synchronized with a sharded parameter server. Model checkpoints are occasionally written to the shared filesystem from the trainers.



BigGraph:

- **Pros:**

- Good powerful engine with high CPU-utilization
- Easy to use
- Implicit weights

- **Cons:**

- Still long for big graphs
- Graph changes -> we need to re-train it each time
- Still some hyper-parameters

- **Personal:**

- tried but got nothing =(
- required ~12h on 63 cores for OK graph
- Restored friends connections really good, not as good in PYMK

Unclassical embeddings

Summary

- Not only MF
- Different embeddings fits different purposes and restrictions
- Based on nature of your product (graph-, session- based)
- It doesn't mean that different embeddings replaces one each other – combine it!
- There are still a lot of tricks to combine one with each other:
 - Similar to last positives
 - Score, is_from_source, num in similar top for each embedding in features