

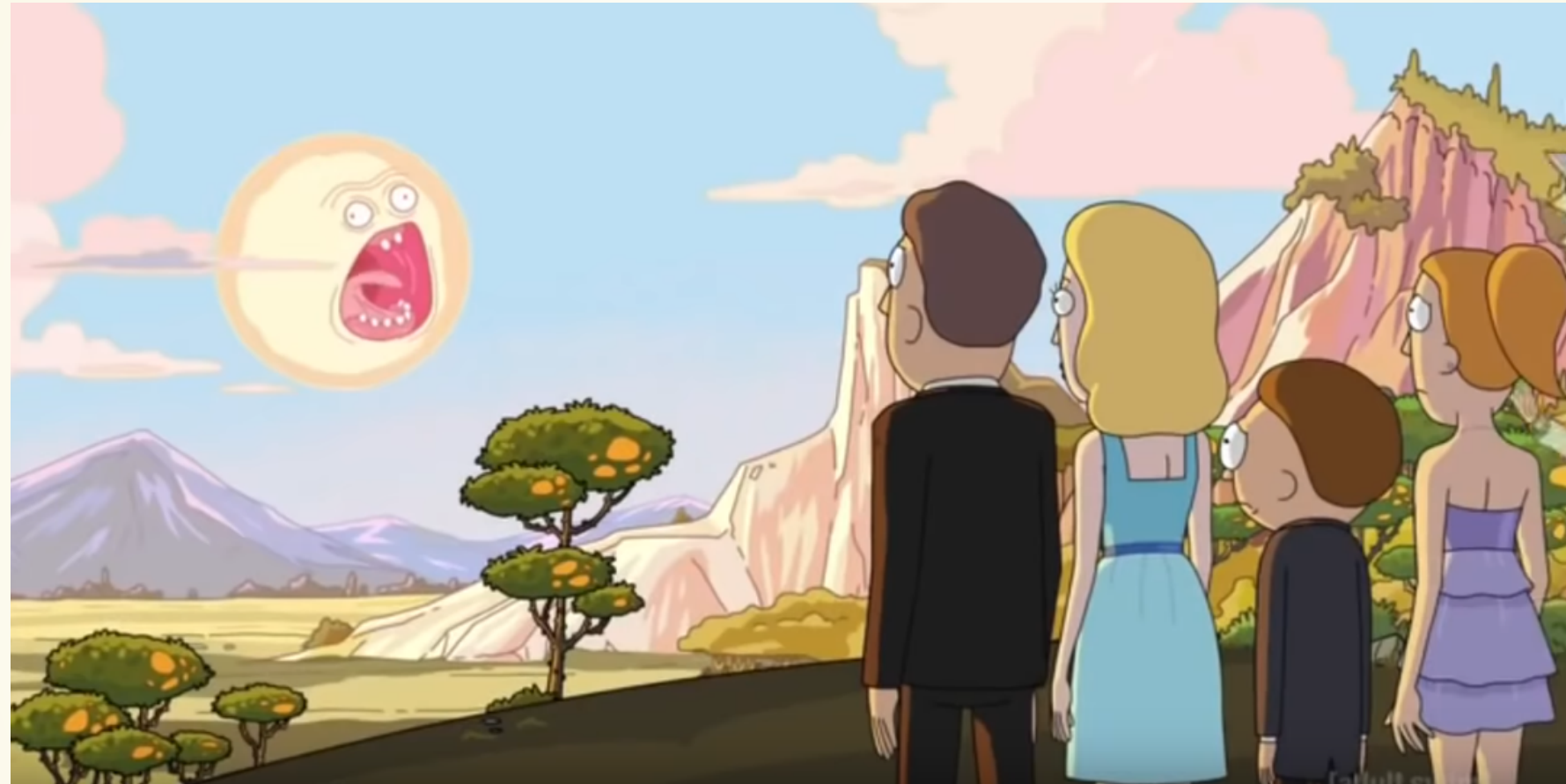


REINFORCEMENT LEARNING

IN SEARCH

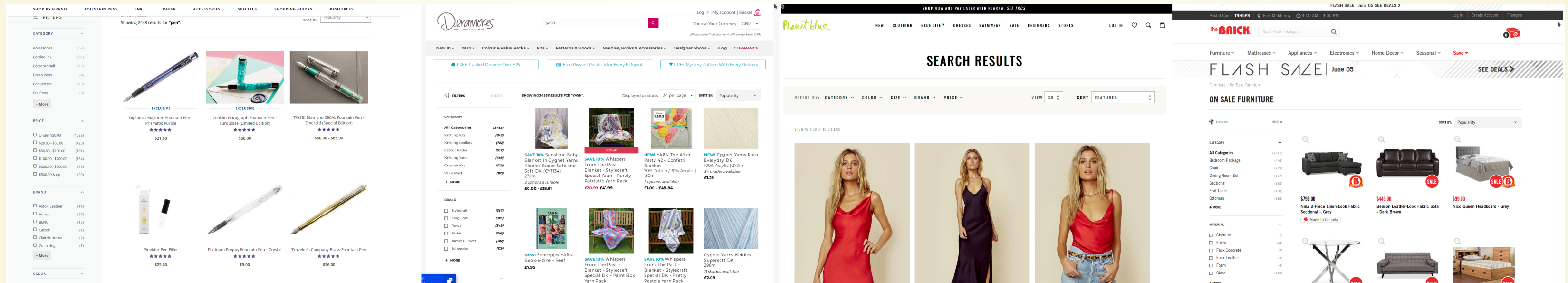
Roman Grebennikov | Findify AB | MICES 2021

ABOUT ME



- Doing ML for 12 years
- Quant trading → Credit risk → eCommerce
- Relevancy & ranking @ Findify

ABOUT FINDIFY



- white-label eCommerce SaaS search
- 20M searchable products
- 50M MAU

AGENDA

- Personalization
- Doing it wrong
- Doing it right *
- Methods and solutions

RANKING: CURRENT STATE

- Search: ElasticSearch/SOLR, TF/IDF, BM25
- Listings: sort by popularity
- Recommendations: CF/SVD/ASL

BM25: an intuitive view

Repetitions of query words \rightarrow good

Common words less important

$$\log \frac{P(D | R=1)}{P(D | R=0)} \approx \sum_w \left(\frac{d_w(1+k)}{d_w + k((1-b) + \frac{b \cdot dl}{avg. dl})} \cdot \log \frac{N - N_w + \frac{1}{2}}{N_w + \frac{1}{2}} \right)$$

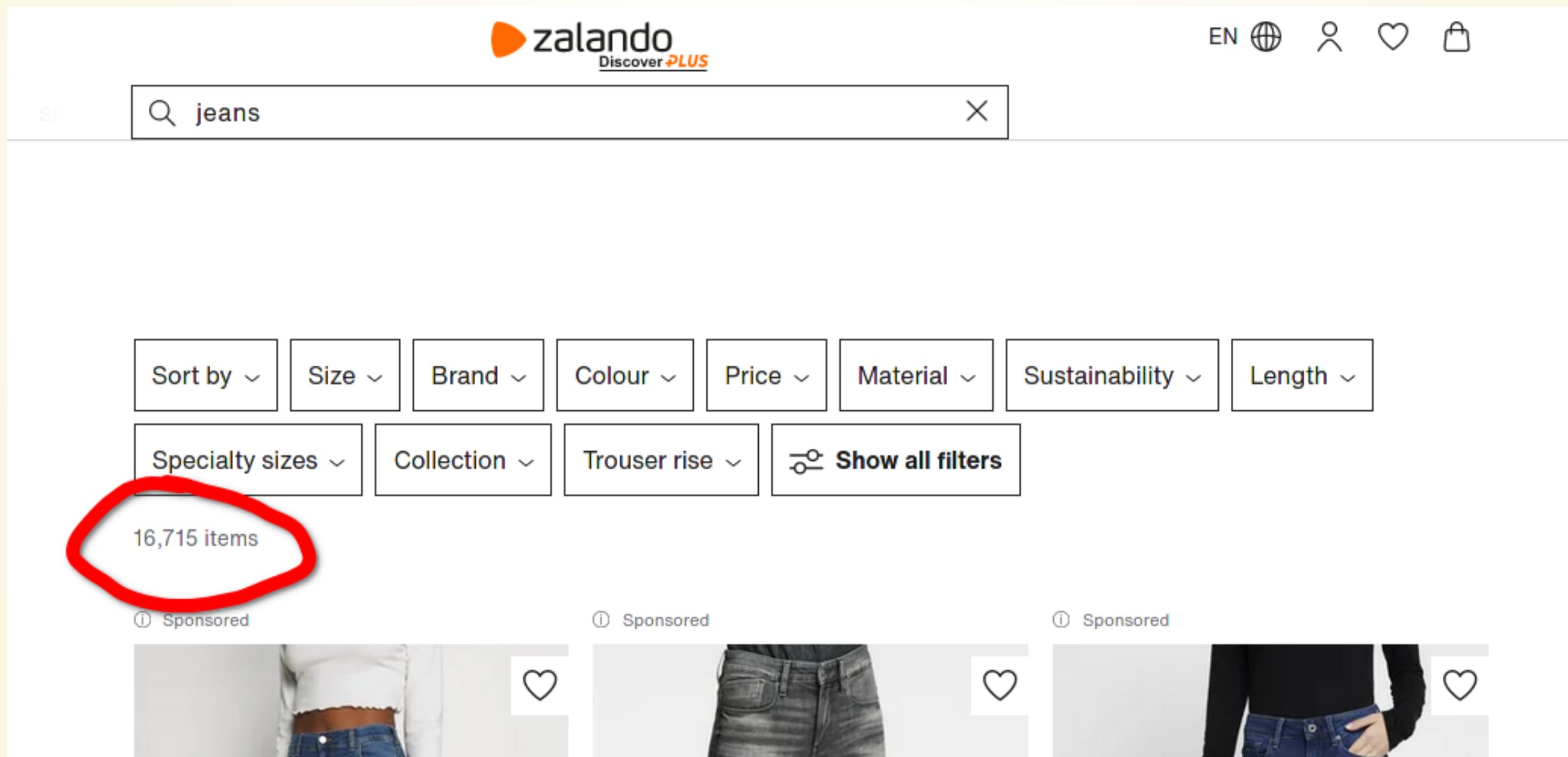
More words in common with the query \rightarrow good

Repetitions less important than different query words

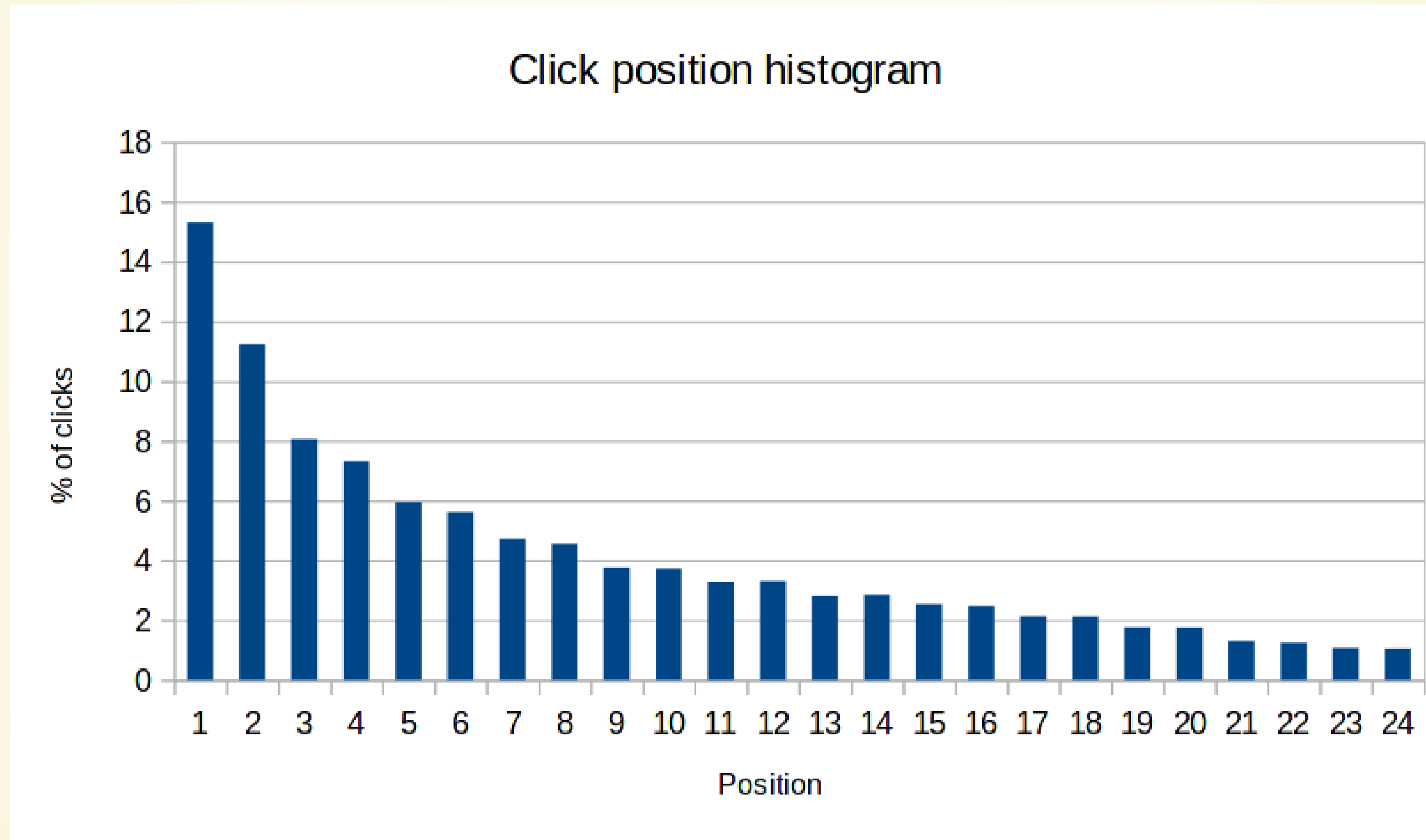
But more important if document is relatively long (wrt. average)

BM25

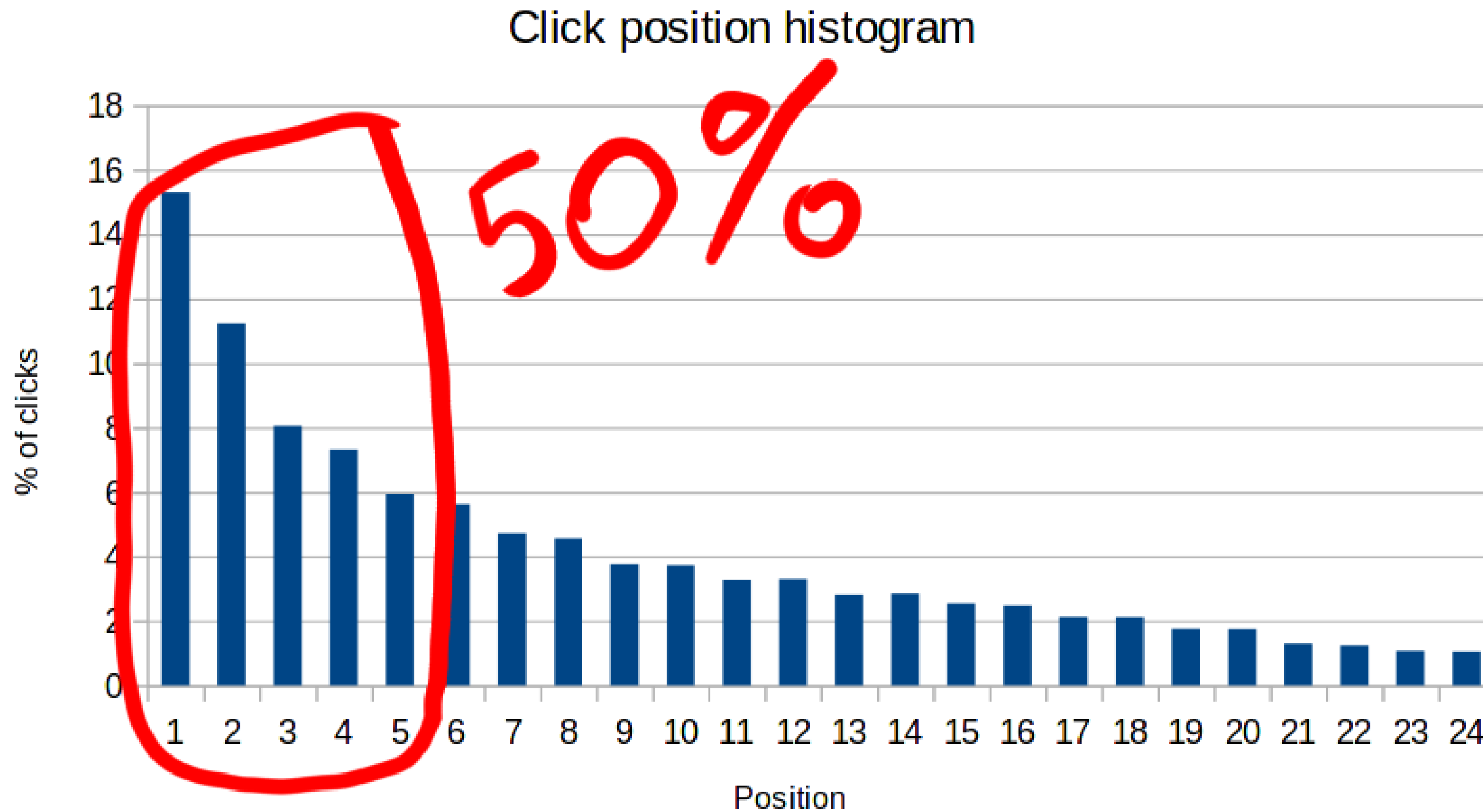
type "cat" - get cat



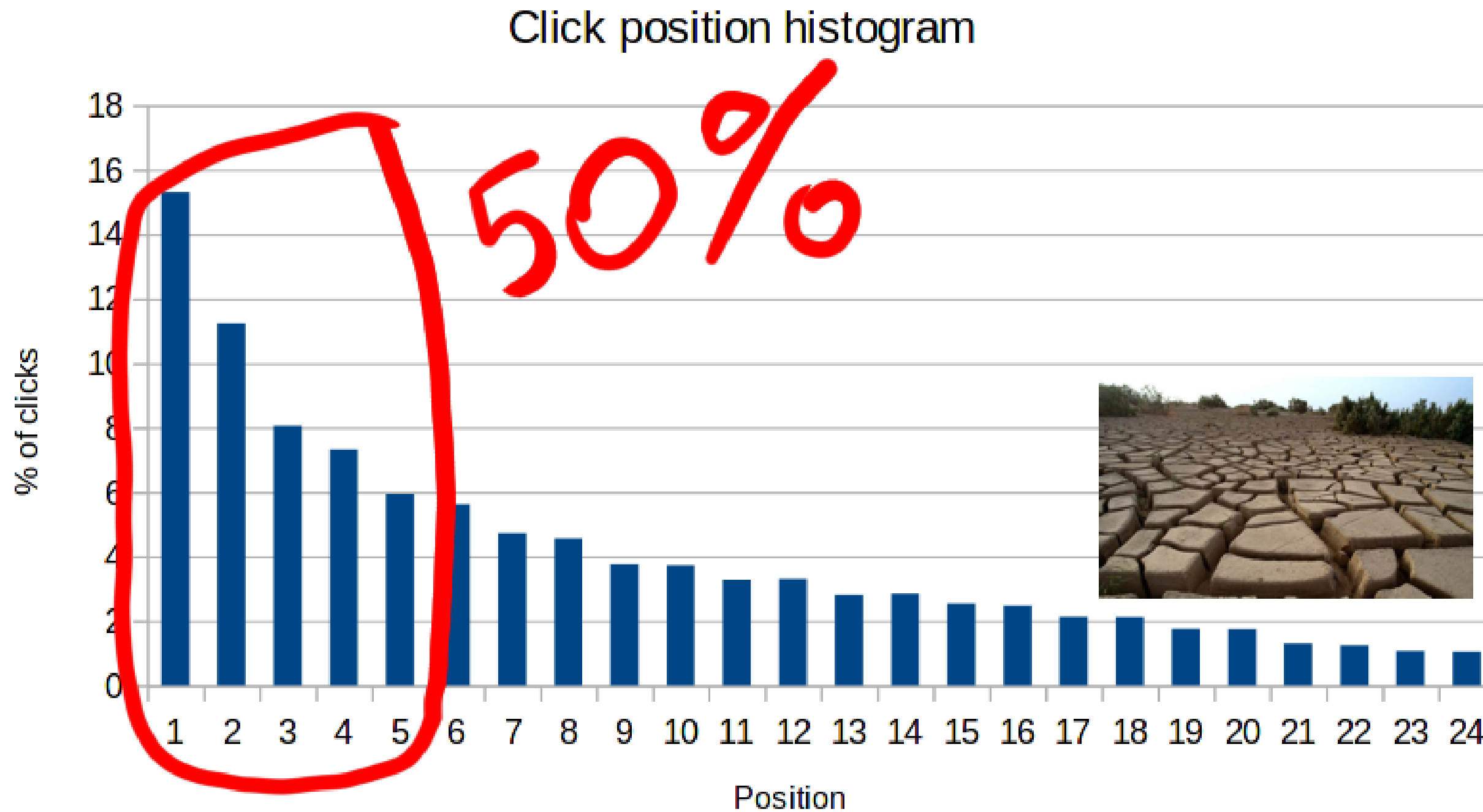
CUSTOMER FOCUS LOST



CUSTOMER FOCUS LOST



CUSTOMER FOCUS LOST



HOW TO LOOSE A CUSTOMER

- only first items are important
- irrelevant top-N = lost customer

BM25 WITH BENEFITS

Idea: rank by relevancy AND popularity

$$\text{score} = \text{bm25} * \log(1 + \text{clicks})$$

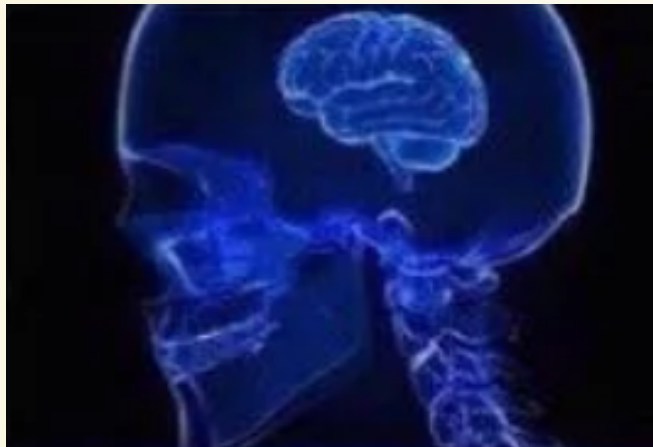
Product	BM25	# of clicks	Score
p1	1.7	20	2.24
p2	1.6	350	4.07
p3	0.7	600	1.94

BM25 WITH BENEFITS

Product	BM25	# of clicks	Score
p2	1.6	350	4.07
p1	1.7	20	2.24
p3	0.7	600	1.94



BM25 WITH BENEFITS



$$\text{score} = \text{bm25}$$



$$\text{score} = \text{bm25} * \log(1 + \text{clicks})$$



$$\text{score} = c_1 + c_2 * \text{bm25} * \log(c_3 + \text{clicks})$$

A BETTER RANKING FORMULA?

- why only clicks?
 - let's add pageviews, purchases and carts!
- let's add windows for 1-7-14 days!
- more constants to tune!



A BETTER RANKING FORMULA?

$$\begin{aligned}
 \mathcal{L}_{SM} = & -\frac{1}{2}\partial_\nu g_\mu^a \partial_\nu g_\mu^a - g_s f^{abc} \partial_\mu g_\nu^a g_\mu^b g_\nu^c - \frac{1}{4}g_s^2 f^{abc} f^{ade} g_\mu^b g_\nu^c g_\mu^d g_\nu^e - \partial_\nu W_\mu^+ \partial_\nu W_\mu^- - \\
 & M^2 W_\mu^+ W_\mu^- - \frac{1}{2}\partial_\nu Z_\mu^0 \partial_\nu Z_\mu^0 - \frac{1}{2c_w^2} M^2 Z_\mu^0 Z_\mu^0 - \frac{1}{2}\partial_\mu A_\nu \partial_\mu A_\nu - ig c_w (\partial_\nu Z_\mu^0 (W_\mu^+ W_\nu^- - W_\nu^+ W_\mu^-) - \\
 & Z_\nu^0 (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) + Z_\mu^0 (W_\nu^+ \partial_\nu W_\mu^- - W_\nu^- \partial_\nu W_\mu^+)) - ig s_w (\partial_\nu A_\mu (W_\mu^+ W_\nu^- - \\
 & W_\nu^+ W_\mu^-) - A_\nu (W_\mu^+ \partial_\nu W_\mu^- - W_\mu^- \partial_\nu W_\mu^+) + A_\mu (W_\nu^+ \partial_\nu W_\mu^- - W_\nu^- \partial_\nu W_\mu^+)) - \\
 & \frac{1}{2}g^2 W_\mu^+ W_\mu^- W_\nu^+ W_\nu^- + \frac{1}{2}g^2 W_\mu^+ W_\nu^- W_\mu^+ W_\nu^- + g^2 c_w^2 (Z_\mu^0 W_\mu^+ Z_\nu^0 W_\nu^- - Z_\mu^0 Z_\mu^0 W_\nu^+ W_\nu^-) + \\
 & g^2 s_w^2 (A_\mu W_\mu^+ A_\nu W_\nu^- - A_\mu A_\mu W_\nu^+ W_\nu^-) + g^2 s_w c_w (A_\mu Z_\nu^0 (W_\mu^+ W_\nu^- - W_\nu^+ W_\mu^-) - \\
 & 2A_\mu Z_\mu^0 W_\nu^+ W_\nu^-) - \frac{1}{2}\partial_\mu H \partial_\mu H - 2M^2 \alpha_h H^2 - \partial_\mu \phi^+ \partial_\mu \phi^- - \frac{1}{2}\partial_\mu \phi^0 \partial_\mu \phi^0 - \\
 & \beta_h \left(\frac{2M^2}{g^2} + \frac{2M}{g} H + \frac{1}{2}(H^2 + \phi^0 \phi^0 + 2\phi^+ \phi^-) \right) + \frac{2M^4}{g^2} \alpha_h - g \alpha_h M (H^3 + H \phi^0 \phi^0 + 2H \phi^+ \phi^-) - \\
 & \frac{1}{8}g^2 \alpha_h (H^4 + (\phi^0)^4 + 4(\phi^+ \phi^-)^2 + 4(\phi^0)^2 \phi^+ \phi^- + 4H^2 \phi^+ \phi^- + 2(\phi^0)^2 H^2) - g M W_\mu^+ W_\mu^- H - \\
 & \frac{1}{2}g \frac{M}{c_w^2} Z_\mu^0 Z_\mu^0 H - \frac{1}{2}ig (W_\mu^+ (\phi^0 \partial_\mu \phi^- - \phi^- \partial_\mu \phi^0) - W_\mu^- (\phi^0 \partial_\mu \phi^+ - \phi^+ \partial_\mu \phi^0)) + \\
 & \frac{1}{2}g (W_\mu^+ (H \partial_\mu \phi^- - \phi^- \partial_\mu H) + W_\mu^- (H \partial_\mu \phi^+ - \phi^+ \partial_\mu H)) + \frac{1}{2}g \frac{1}{c_w} (Z_\mu^0 (H \partial_\mu \phi^0 - \phi^0 \partial_\mu H) + \\
 & M (\frac{1}{c_w} Z_\mu^0 \partial_\mu \phi^0 + W_\mu^+ \partial_\mu \phi^- + W_\mu^- \partial_\mu \phi^+) - ig \frac{s_w}{c_w} M Z_\mu^0 (W_\mu^+ \phi^- - W_\mu^- \phi^+) + ig s_w M A_\mu (W_\mu^+ \phi^- - \\
 & W_\mu^- \phi^+) - ig \frac{1-2c_w^2}{2c_w} Z_\mu^0 (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) + ig s_w A_\mu (\phi^+ \partial_\mu \phi^- - \phi^- \partial_\mu \phi^+) - \\
 & \frac{1}{4}g^2 W_\mu^+ W_\mu^- (H^2 + (\phi^0)^2 + 2\phi^+ \phi^-) - \frac{1}{8}g^2 \frac{1}{c_w^2} Z_\mu^0 Z_\mu^0 (H^2 + (\phi^0)^2 + 2(2s_w^2 - 1)^2 \phi^+ \phi^-) - \\
 & \frac{1}{2}g^2 \frac{s_w^2}{c_w} Z_\mu^0 \phi^0 (W_\mu^+ \phi^- + W_\mu^- \phi^+) - \frac{1}{2}ig^2 \frac{s_w^2}{c_w} Z_\mu^0 H (W_\mu^+ \phi^- - W_\mu^- \phi^+) + \frac{1}{2}g^2 s_w A_\mu \phi^0 (W_\mu^+ \phi^- + \\
 & W_\mu^- \phi^+) + \frac{1}{2}ig^2 s_w A_\mu H (W_\mu^+ \phi^- - W_\mu^- \phi^+) - g^2 \frac{s_w}{c_w} (2c_w^2 - 1) Z_\mu^0 A_\mu \phi^+ \phi^- - g^2 s_w^2 A_\mu A_\mu \phi^+ \phi^- + \\
 & \frac{1}{2}ig s_\lambda \lambda_{ij}^a (\bar{q}_i^\sigma \gamma^\mu q_j^\sigma) g_\mu^a - \bar{e}^\lambda (\gamma \partial + m_e^\lambda) e^\lambda - \bar{\nu}^\lambda (\gamma \partial + m_\nu^\lambda) \nu^\lambda - \bar{u}_j^\lambda (\gamma \partial + m_u^\lambda) u_j^\lambda - \bar{d}_j^\lambda (\gamma \partial + m_d^\lambda) d_j^\lambda + \\
 & ig s_w A_\mu \left(-(\bar{e}^\lambda \gamma^\mu e^\lambda) + \frac{2}{3}(\bar{u}_j^\lambda \gamma^\mu u_j^\lambda) - \frac{1}{3}(\bar{d}_j^\lambda \gamma^\mu d_j^\lambda) \right) + \frac{ig}{4c_w} Z_\mu^0 \{ (\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{e}^\lambda \gamma^\mu (4s_w^2 - \\
 & 1 - \gamma^5) e^\lambda) + (\bar{d}_j^\lambda \gamma^\mu (\frac{4}{3}s_w^2 - 1 - \gamma^5) d_j^\lambda) + (\bar{u}_j^\lambda \gamma^\mu (1 - \frac{8}{3}s_w^2 + \gamma^5) u_j^\lambda) \} + \\
 & \frac{ig}{2\sqrt{2}} W_\mu^+ \left((\bar{\nu}^\lambda \gamma^\mu (1 + \gamma^5) U^{lep}_{\lambda\kappa} e^\kappa) + (\bar{u}_j^\lambda \gamma^\mu (1 + \gamma^5) C_{\lambda\kappa} d_j^\kappa) \right) + \\
 & \frac{ig}{2\sqrt{2}} W_\mu^- \left((\bar{e}^\kappa U^{lep\dagger}_{\kappa\lambda} \gamma^\mu (1 + \gamma^5) \nu^\lambda) + (\bar{d}_j^\kappa C_{\kappa\lambda}^\dagger \gamma^\mu (1 + \gamma^5) u_j^\lambda) \right) + \\
 & \frac{ig}{2M\sqrt{2}} \phi^+ \left(-m_e^\kappa (\bar{\nu}^\lambda U^{lep}_{\lambda\kappa} (1 - \gamma^5) e^\kappa) + m_\nu^\lambda (\bar{\nu}^\lambda U^{lep}_{\lambda\kappa} (1 + \gamma^5) e^\kappa) + \right. \\
 & \left. \frac{ig}{2M\sqrt{2}} \phi^- \left(m_e^\lambda (\bar{e}^\lambda U^{lep\dagger}_{\lambda\kappa} (1 + \gamma^5) \nu^\kappa) - m_\nu^\kappa (\bar{e}^\lambda U^{lep\dagger}_{\lambda\kappa} (1 - \gamma^5) \nu^\kappa) - \frac{g}{2} \frac{m_\nu^\lambda}{M} H (\bar{\nu}^\lambda \nu^\lambda) - \right. \right. \\
 & \left. \frac{g}{2} \frac{m_\lambda^\lambda}{M} H (\bar{e}^\lambda e^\lambda) + \frac{ig}{2} \frac{m_\lambda^\lambda}{M} \phi^0 (\bar{\nu}^\lambda \gamma^5 \nu^\lambda) - \frac{ig}{2} \frac{m_\lambda^\lambda}{M} \phi^0 (\bar{e}^\lambda \gamma^5 e^\lambda) - \frac{1}{4} \bar{\nu}_\lambda M_{\lambda\kappa}^R (1 - \gamma_5) \hat{\nu}_\kappa - \right. \\
 & \left. \frac{1}{4} \bar{\nu}_\lambda M_{\lambda\kappa}^R (1 - \gamma_5) \hat{\nu}_\kappa + \frac{ig}{2M\sqrt{2}} \phi^+ \left(-m_d^\kappa (\bar{u}_j^\lambda C_{\lambda\kappa} (1 - \gamma^5) d_j^\kappa) + m_u^\lambda (\bar{u}_j^\lambda C_{\lambda\kappa} (1 + \gamma^5) d_j^\kappa) \right) + \right. \\
 & \left. \frac{ig}{2M\sqrt{2}} \phi^- \left(m_d^\lambda (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 + \gamma^5) u_j^\kappa) - m_u^\kappa (\bar{d}_j^\lambda C_{\lambda\kappa}^\dagger (1 - \gamma^5) u_j^\kappa) - \frac{g}{2} \frac{m_\lambda^\lambda}{M} H (\bar{u}_j^\lambda u_j^\lambda) - \frac{g}{2} \frac{m_\lambda^\lambda}{M} H (\bar{d}_j^\lambda d_j^\lambda) + \right. \right. \\
 & \left. \left. \frac{ig}{2} \frac{m_\lambda^\lambda}{M} \phi^0 (\bar{u}_j^\lambda \gamma^5 u_j^\lambda) - \frac{ig}{2} \frac{m_\lambda^\lambda}{M} \phi^0 (\bar{d}_j^\lambda \gamma^5 d_j^\lambda) \right) \right.
 \end{aligned}$$

probably won't work out of the box -> needs training!



LEARN-TO-RANK

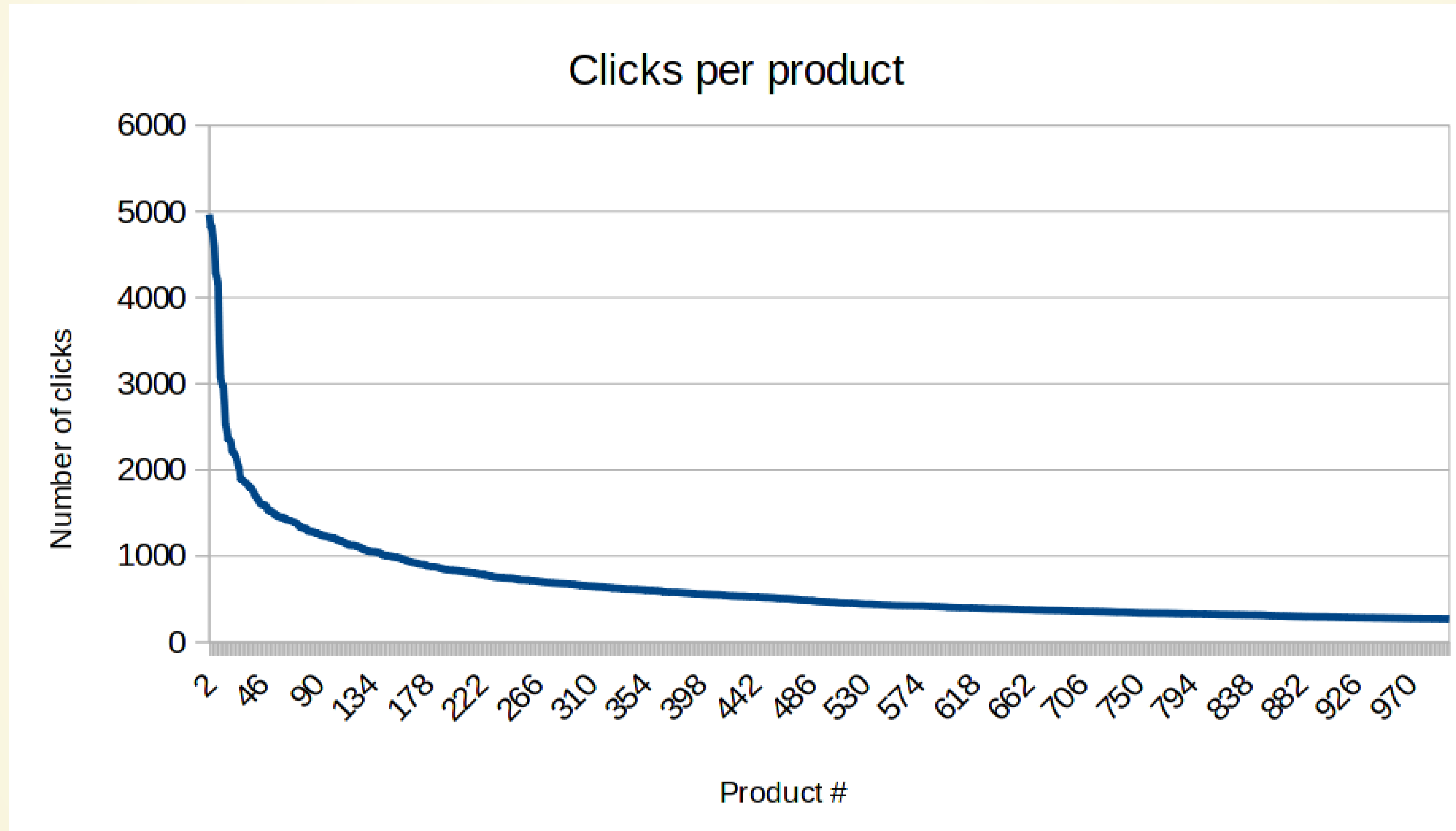
**ME WITH
ELASTICSEARCH**

A. WALLIN

LTR: GOAL TO OPTIMIZE

Product	BM25	# of clicks	Interaction
p1	1.7	100	-
p2	0.9	20	click
p3	0.7	600	-

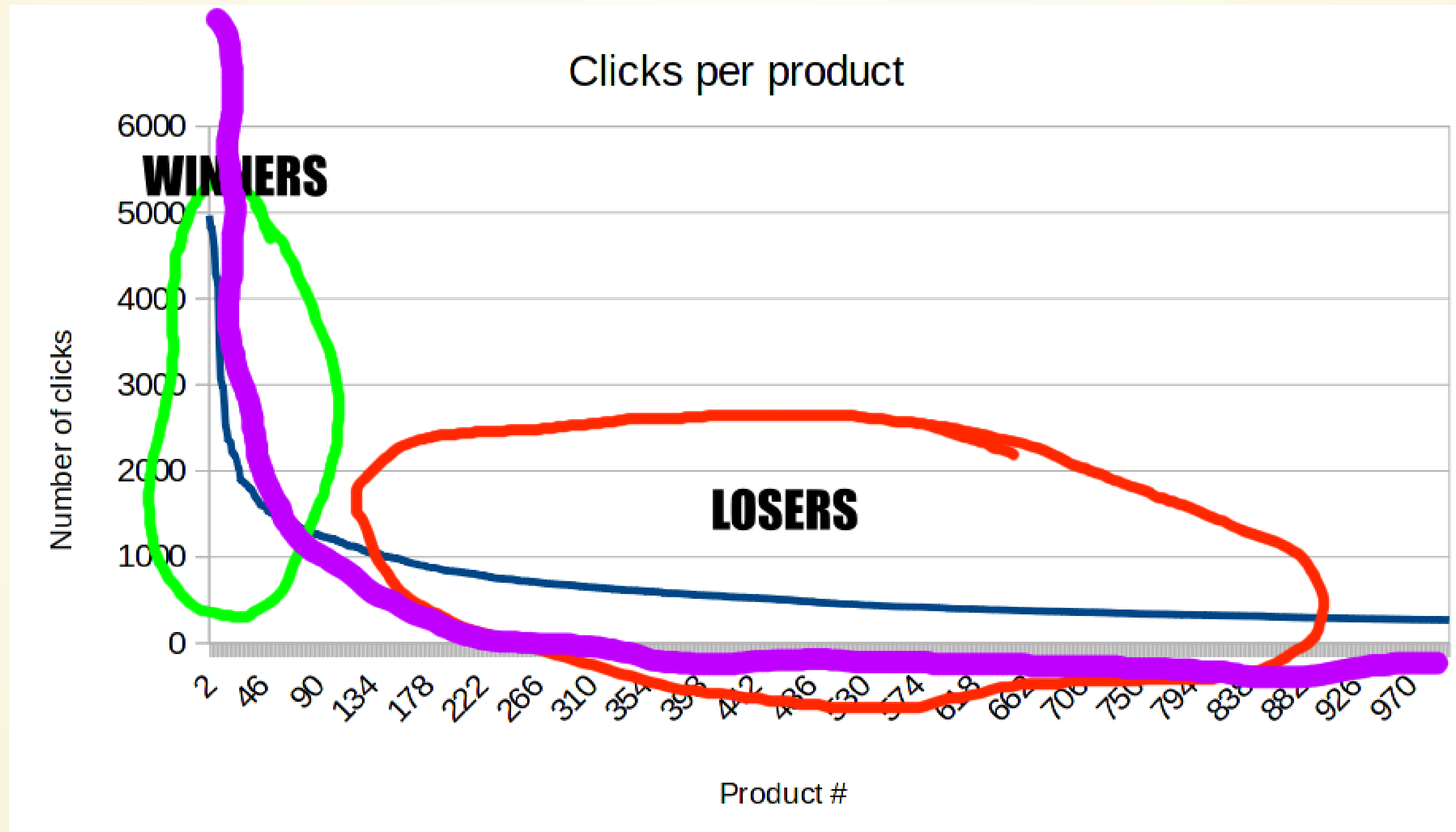
FEEDBACK LOOP



FEEDBACK LOOP OF BAD SEARCH

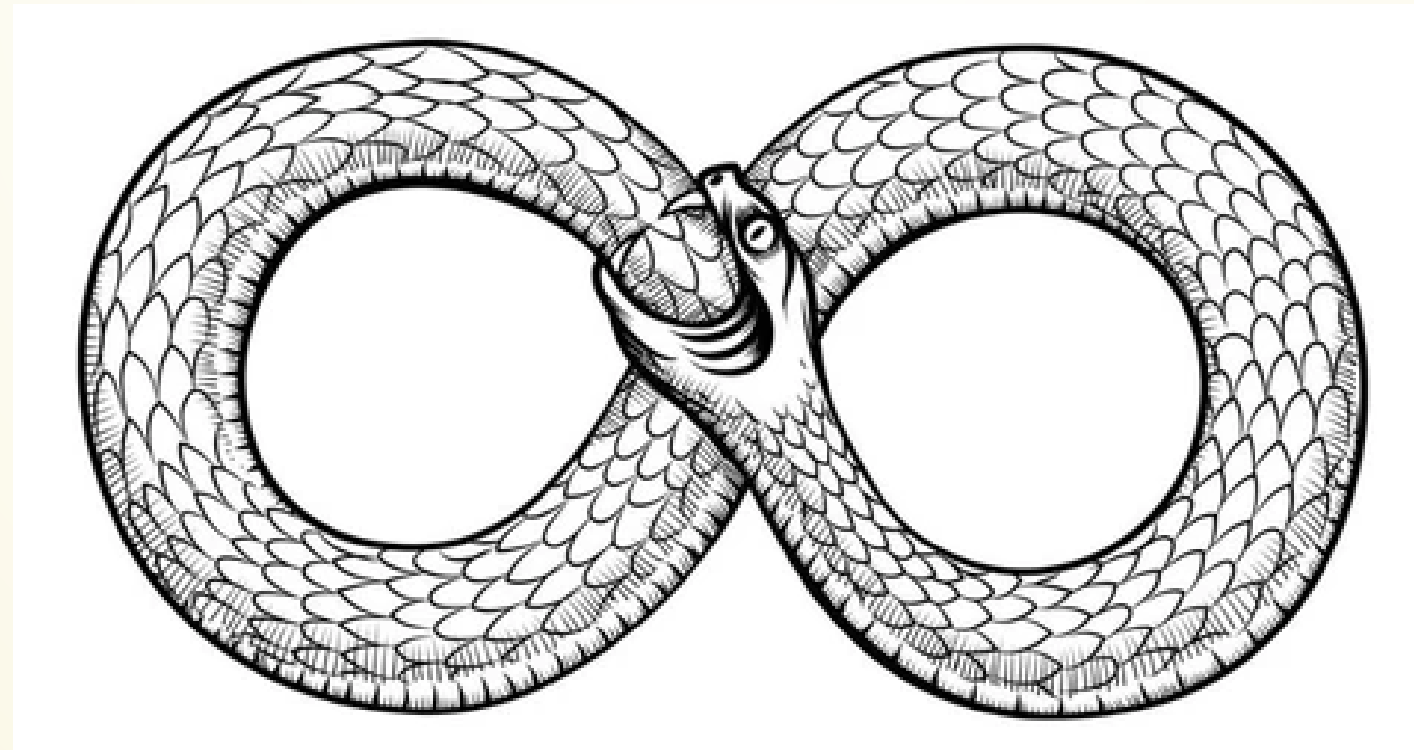


FEEDBACK LOOP OF BAD SEARCH

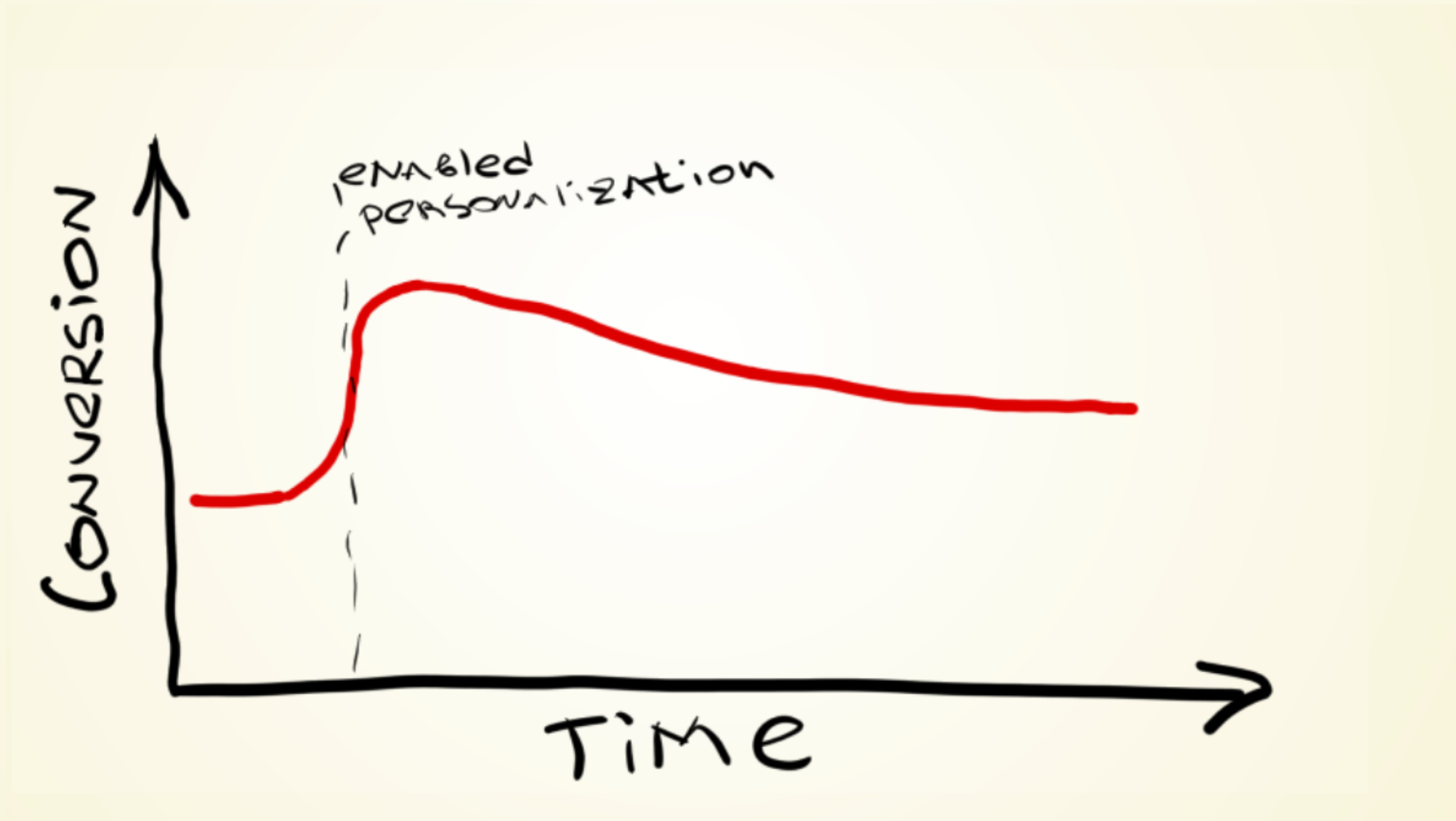


MODEL BIAS

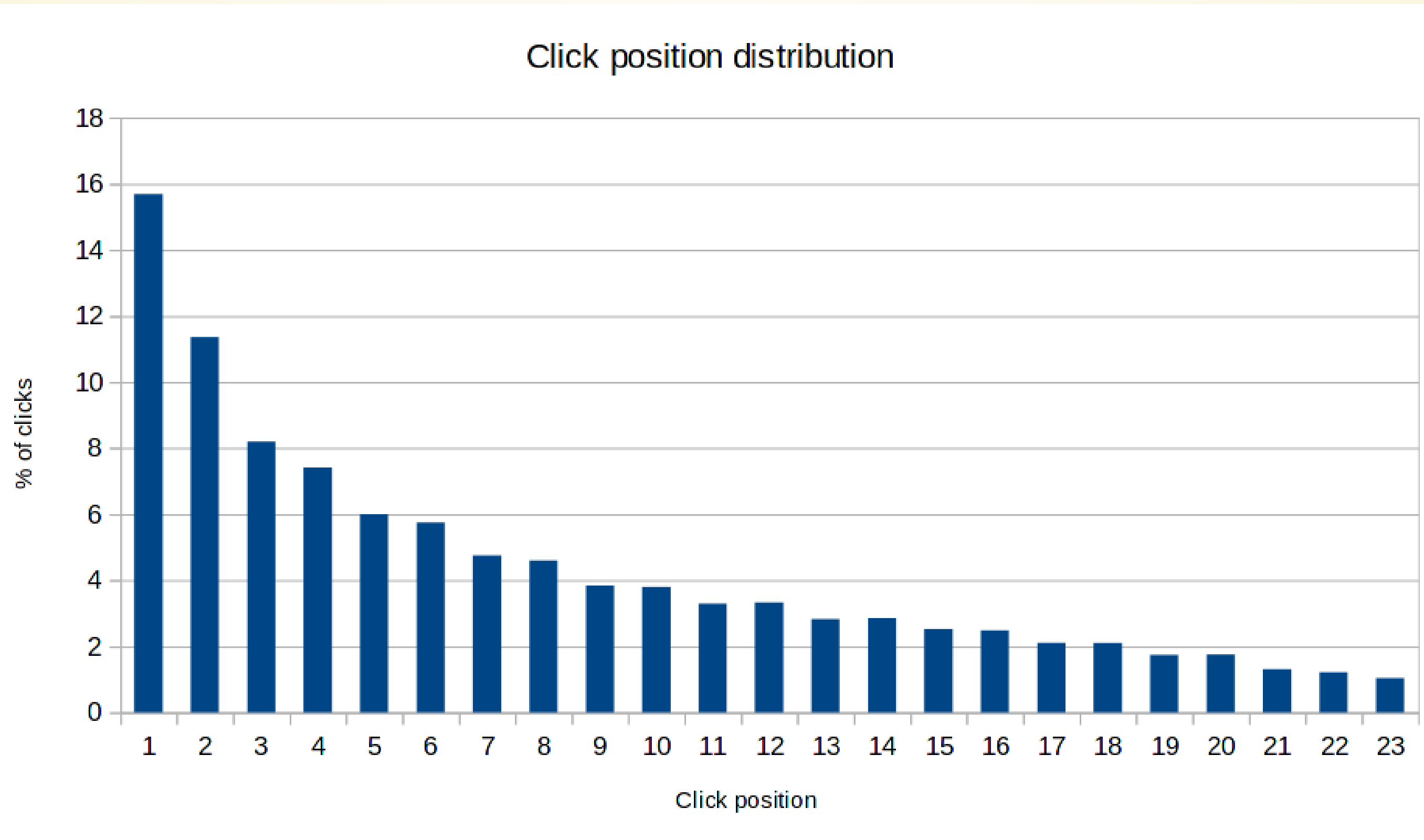
- Produce ranking based on past clicks
- Collect clicks influenced by previous ranking
- Retrain model on clicks
- Repeat



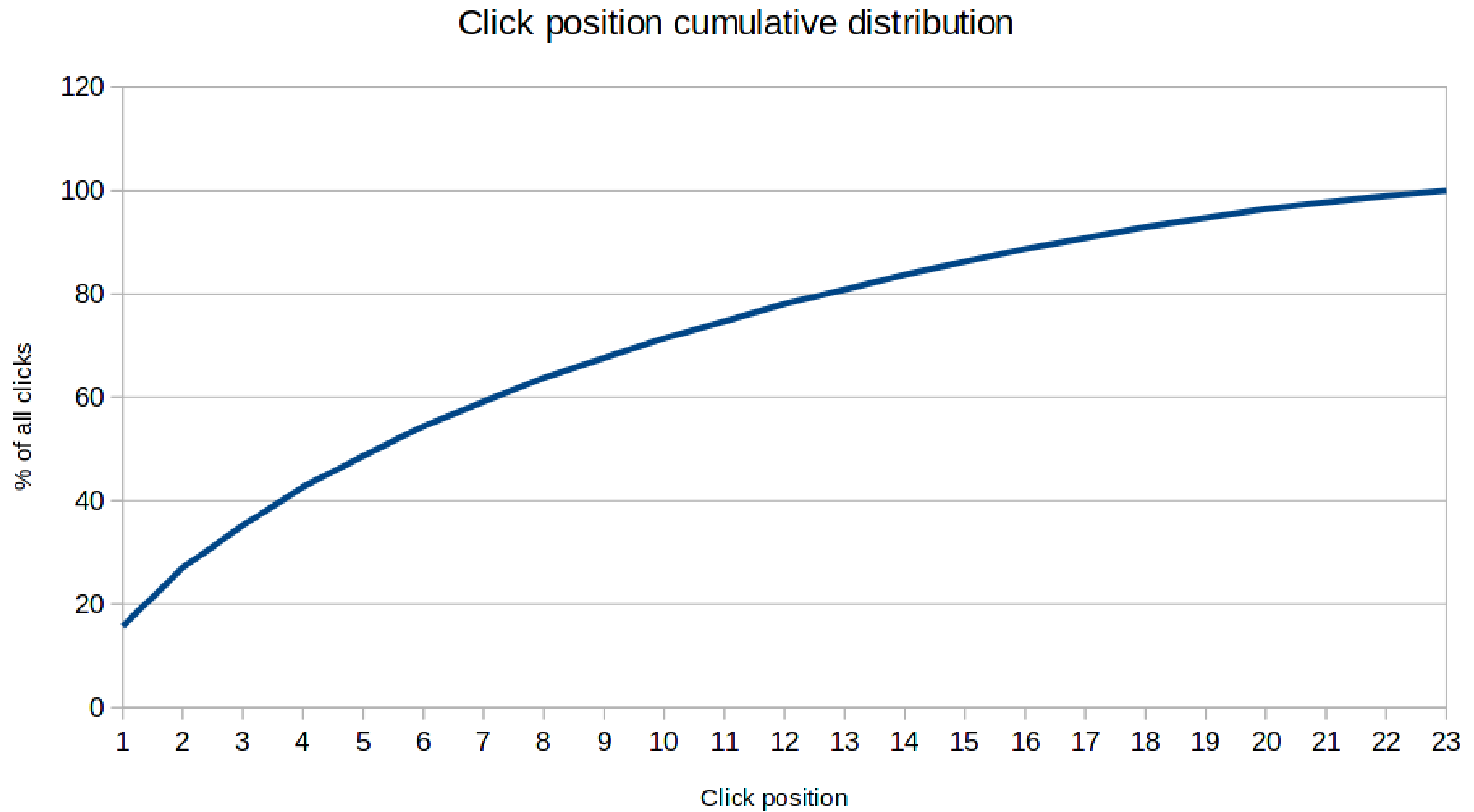
TRUE PRODUCTION STORY



POSITION BIAS

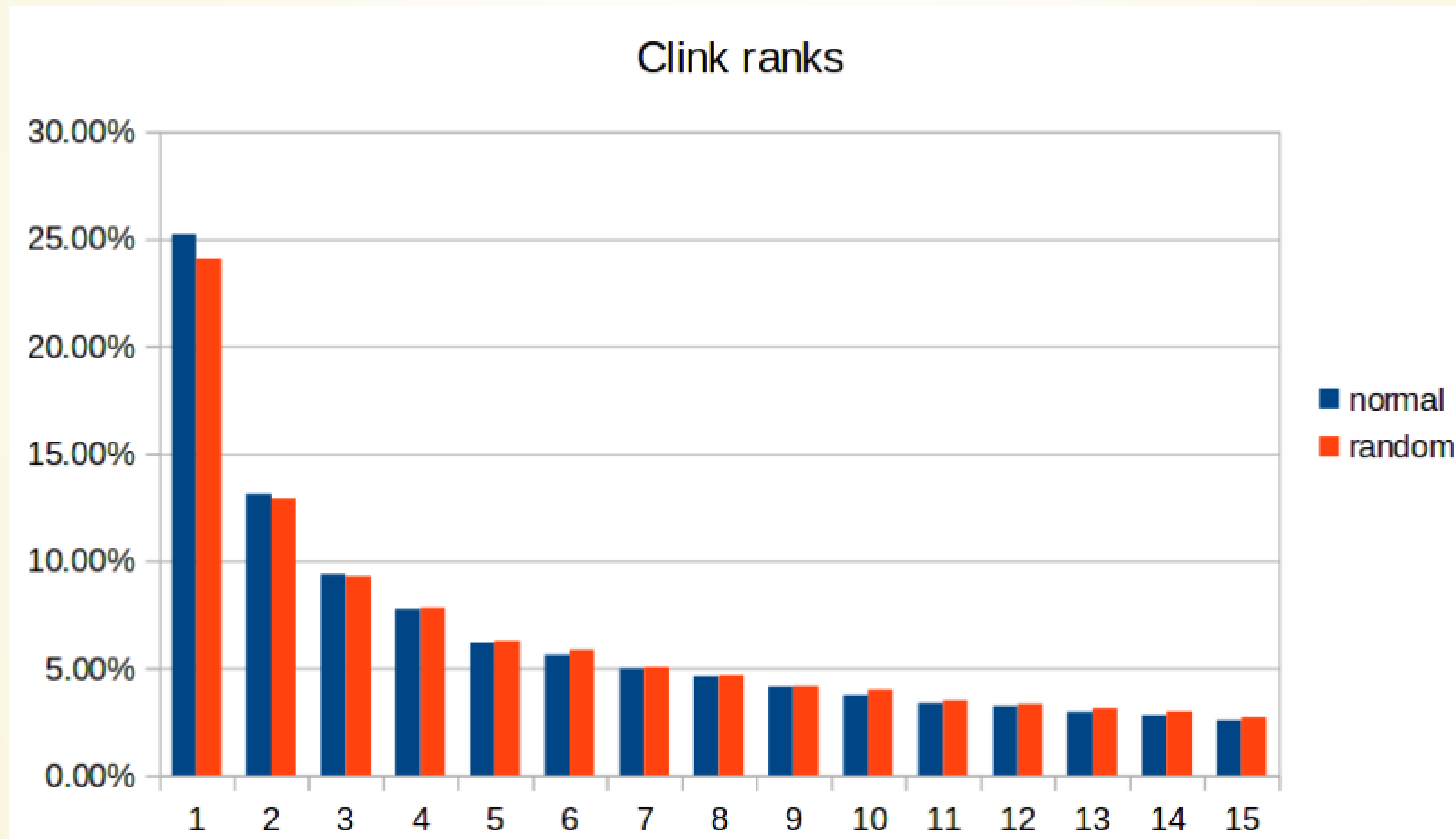


POSITION BIAS



POSITION BIAS: EXPERIMENT

default ordering vs random shuffling a/b test





ESTIMATE THE BIAS?

observed relevancy = bias + true relevancy

1. Estimate the bias
2. Subtract the bias

...

3. ~~PROFIT~~ true relevancy

IPW: INVERSE PROPENSITY WEIGHTING

Wang X. - Position bias estimation for unbiased learning to rank in personal search

- Query-level Propensity:
 - downsample queries where $\text{avg}(\text{click_rank})$ is too high
- Document-level Propensity:
 - lower weight for top clicks

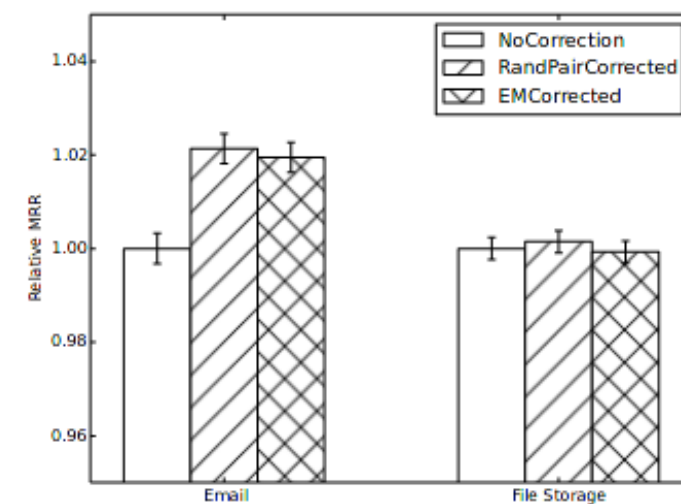


Figure 3: Ranking effectiveness comparison of pairwise approaches with different position bias correction methods.



DOCUMENT LEVEL PROPENSITY

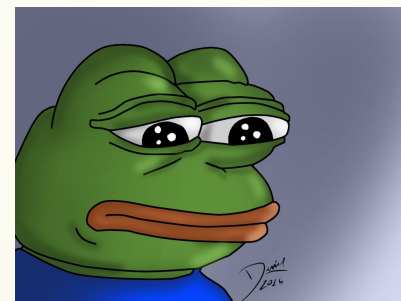
Query: socks

- "Star wars socks" on position 1: 50 clicks
- "Star wars socks" on position 5: 20 clicks
- "Star wars socks" on position 25: 3 clicks

IPW: PROS AND CONS

- Estimation depends on context
- More contexts - more data needed
- Improper estimation may produce another bias
- Need for shuffling

You're not google!



CURSE OF OFFLINE LEARNING

- production training data is biased
- de-biasing is hard

do we really need to stay offline?

REINFORCEMENT LEARNING



- Explore: poke real people, observe reaction
 - if you rank by popularity, noone knows what would happen if ranked by CTR
- Exploit: use it for better results
- Goal: minimize the **regret** / maximize the **reward**

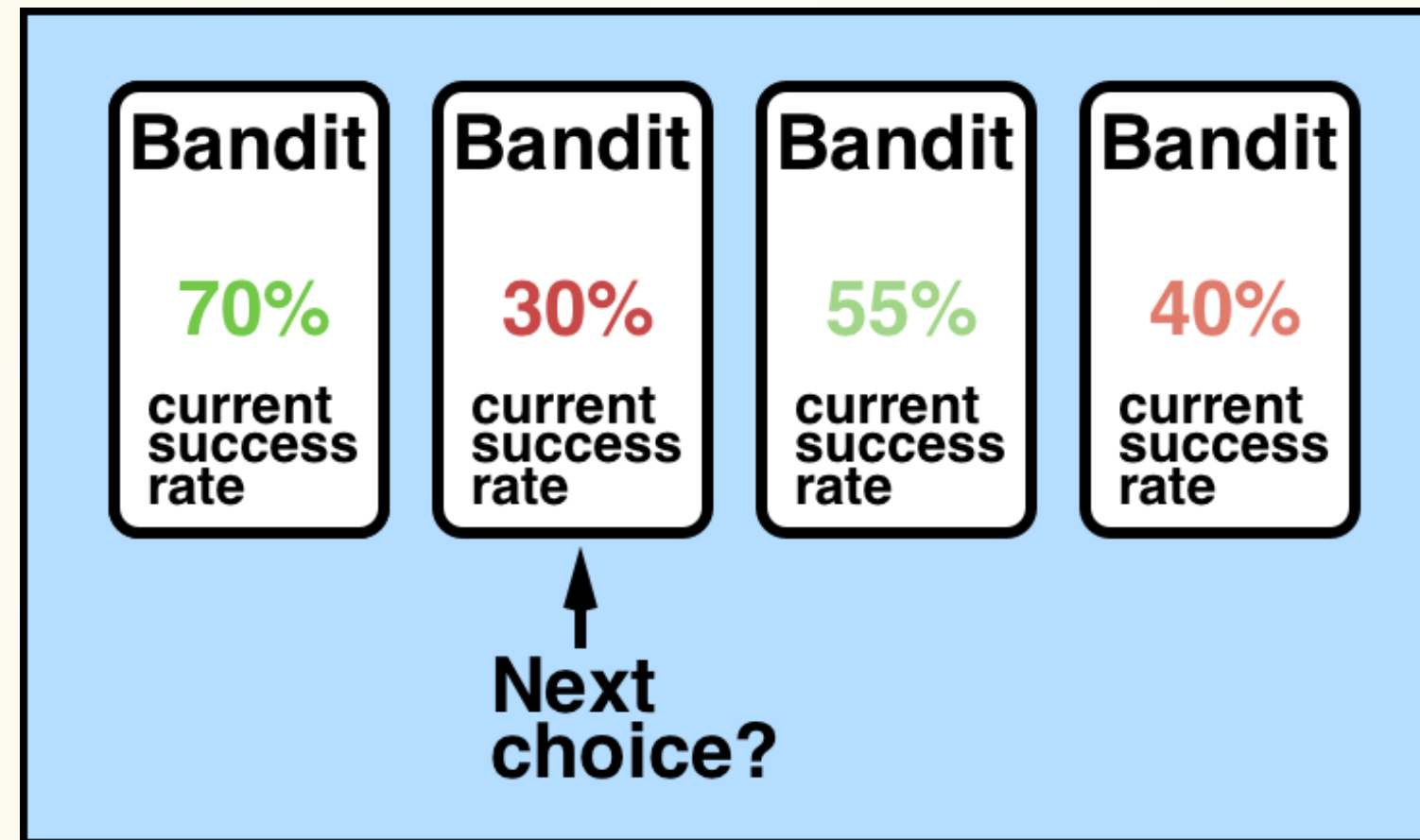
MULTI-ARMED BANDITS



MULTI-ARMED BANDITS

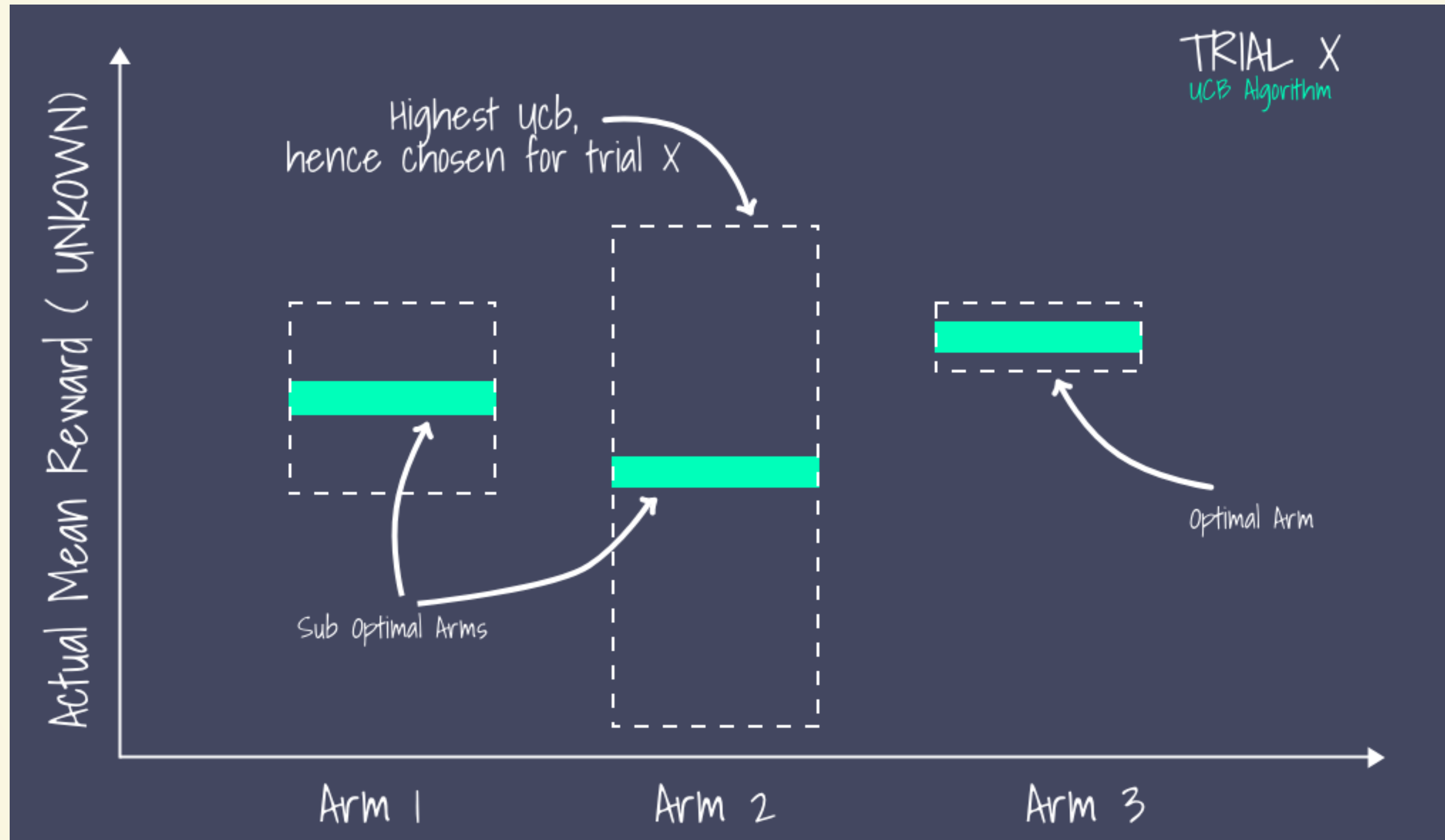
- Slot machines with unknown win probability
- Exploit: maximise the \$\$\$ won
- Explore: which one has highest win chance?

MULTI-ARMED BANDITS



- Greedy: choose arm with max success rate
- ϵ -greedy: $P(\epsilon) = \text{random}$, $P(1-\epsilon) = \text{greedy}$
- LinUCB: choose arm with Upper Confidence Bound

LIN-UCB



LIN-UCB

- Choose the arm with the highest optimistic reward
- Won? UCB goes up, eat cookie



- Lost? UCB goes down
=> Other arm on next iteration

LIN-UCB MEETS RANKING

What is an arm?

- Document?
 - Works great for small datasets like news articles
 - What if item set is large?
- Feature?
 - How rank influences reward?

CASCADE-LIN-UCB

S. Zhong: Cascading Bandits for Large-Scale Recommendation Problems

$$\text{score} = c_1 * f_1 + \dots + c_N * f_N$$

•

- Arm: feature value (# of clicks, purchases)
- Explore: increase-decrease feature weight
 - rank by score, see how it affects reward
- Reward: Increase \$\$\$
- Choose next arm via LinUCB

RL: A/B TEST ON STERIODS

- A/B test when segments are dynamic
- Compare previous iteration with candidate
- Reward goes up = success



CASCADE-LIN-UCB IN PRACTICE

- Reward doesn't depend on score
 - Combine multiple business objectives!
 - 50% CTR + 50% Conversion
 - Click probability?
- Linear model: easy to explain
- Requires a LOT of time to converge
- Linear model: tough with non-linearities

GOING NDCG

H. Oosterhuis: Differentiable Unbiased Online Learning to Rank

- Reward \approx NDCG
- NDCG is not smooth \Rightarrow use pairwise loss

PDGD IN PRACTICE

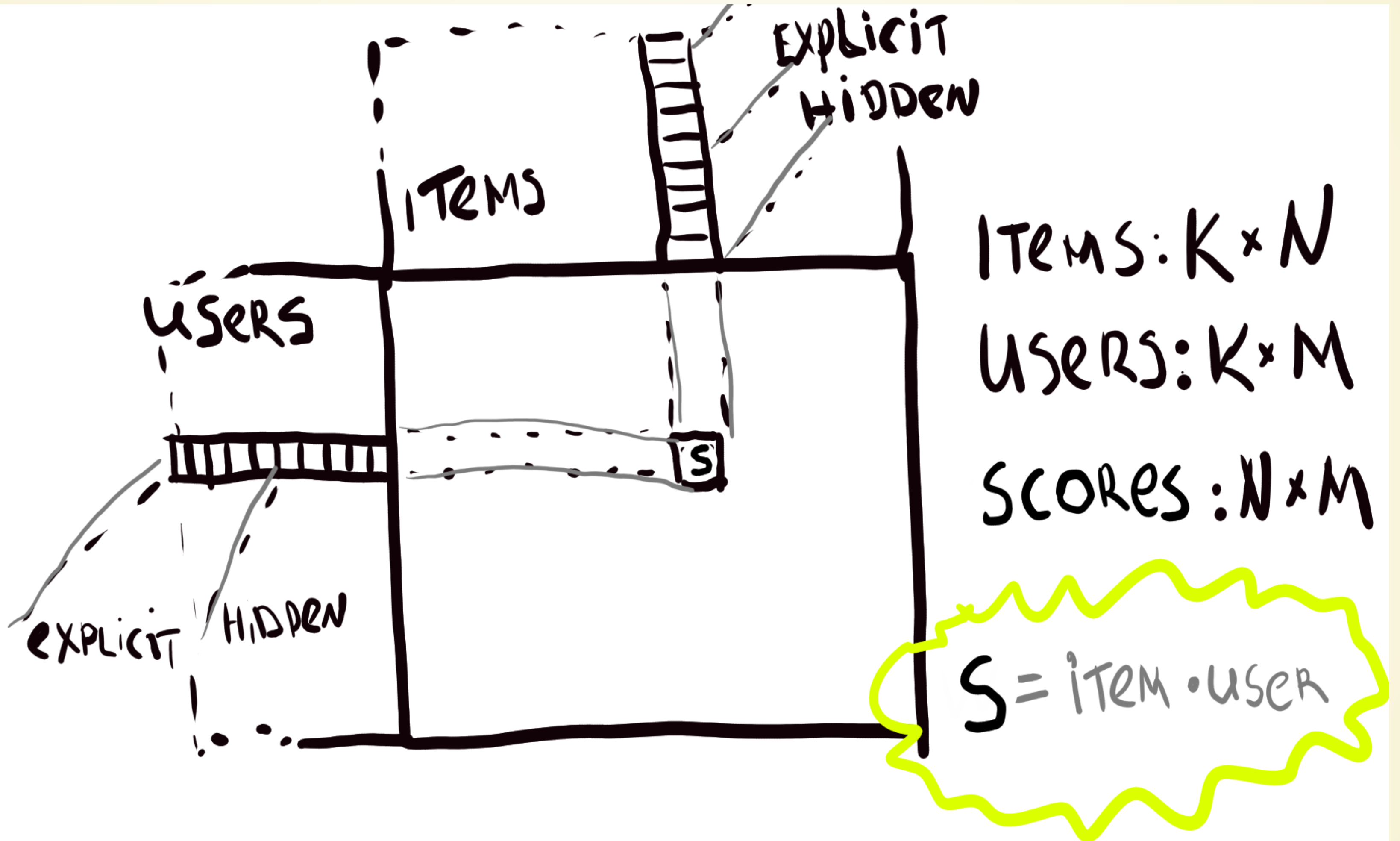
- Much faster to converge
 - weights are updated on each interaction!
- No way to plug whatever reward :(
 - Focused on pairwise loss => better NDCG
- Linear model: easy to explain
- Linear model: tough with non-linearities



COLLABORATIVE RANKING

S. Gupta @ Amazon: CPR: Collaborative Pairwise Ranking for Online List Recommendations

- User/item explicit and latent features
- Weight \sim pairwise loss
- Like MF SVD/ALS for recommendations



CPR IN PRACTICE

- Complex math: tricky to update in realtime
- Focus on pairwise loss
- Good luck with explainability

WELCOME TO THE LEARN-TO-RANK



PLEASE FOLLOW ME

CURRENT STATE OF AFFAIRS



- Heavy NDCG focus
- NDCG is not a business metric!
 - There were cases when NDCG goes up, but CTR & Conversion goes down
 - What if people click higher but less frequent?
- Focusing on AOV/Conv/CTR is hard :(

TOWARDS A SMARTER RANKING

- Which reward is important to you?
- Ensemble of models: combine strengths!

LINKS



- Slides: dfdx.me/slides/reinforcement-learning-in-search
- Me: linkedin.com/in/romangrebennikov/

QUESTIONS?

Derp Learning	100	200	300	400	500
Holywar	100	200	300	400	500
Code of Conduct	100	200	300	400	500
Cargo cult	100	200	300	400	500
TensorFlow	100	200	300	400	500
Haskell	100	200	300	400	500

a slide by @nikitonsky