

Evaluating Search Engines by Modeling the Relationship Between Relevance and Clicks

Ben Carterette & Rosie Jones

Presented by Congcong Li

04 / 01 / 2010

Introduction

- Goal: Evaluate different search engines

Query: Structured SVM

Measure Relevance

Introduction

- How to evaluate search engines?
 - Use relevance judgments
 - Complete relevance judgments not available
- How to predict relevance?
 - Use clicks
 - General problem: Clicks are biased!
- This work:
 - Model the relationship between relevance & clicks**

Evaluating Search Engine

- Discount Cumulative Gain

$$DCG_l = \sum_{i=1}^l gain_i \cdot discount_i$$

$$= rel_1 + \sum_{i=2}^l rel_i \cdot \frac{1}{\log_2 i}$$

$rel_i \in \{a_1, a_2, a_3, a_4, a_5\}$

{Perfect, Excellent, Good, Fair, Bad}

Rank	Relevance	Discount
1	Perfect	1
2	Excellent	1/2
3	Good	1/3
4	Fair	1/4
5	Bad	1/5
6	Bad	1/6

Evaluating DCG from Incomplete Information

- Express DCG as a random variable:

$$DCG_l = X_1 + \sum_{i=2}^l X_i \cdot \frac{1}{\log_2 i} \quad X_i \in \{a_1, a_2, a_3, a_4, a_5\}$$
- Expectation and variance (to estimate confidence interval):

$$E[DCG_l] = E[X_1] + \sum_{i=2}^l \frac{E[X_i]}{\log_2 i}$$

$$Var[DCG_l] = Var[X_1] + \sum_{i=2}^l \frac{E[X_i]}{(\log_2 i)^2} + 2 \sum_{i=2}^l \frac{Cov(X_1, X_i)}{\log_2 i} + 2 \sum_{i < j} \frac{Cov(X_i, X_j)}{\log_2 i \cdot \log_2 j}$$

= 0 under independent assumption
- How to get $p(X_i = a_j)$? **Predict from clicks!**

Modeling Clicks & Relevance

- We want to predict $p(X_i = a_j)$

Joint probability $p(q, X, c)$

$p(q, X_1, X_2, \dots, X_i, c_1, c_2, \dots, c_i)$

↑ query relevance ↑ clickthrough rates

Conditional probability $p(X | q, c)$

Independent assumption

$$p(X | q, c) = \prod_{i=1}^l p(X_i | q, c)$$

query

X_1, c_1

X_2, c_2

X_3, c_3

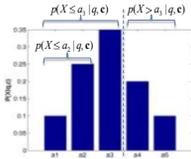
X_4, c_4

X_5, c_5

Modeling Clicks & Relevance

- Model $p(X | q, c) = \prod_{i=1}^l p(X_i | q, c)$
- Ordinal Regression

$$\log \frac{p(X > a_j | q, c)}{p(X \leq a_j | q, c)} = a_j + \beta q + \sum_{i=1}^l \beta_i c_i + \sum_{i < k} \beta_{ik} c_i c_k$$



Scalars
 a_j : one of the relevance levels
 q : the aggregate click rate over all results
 c_i : click times over # of times the list was shown

- Vector Generalized additive model (VGAM)

Comparative Evaluation

- If we only care about whether one ranking function outperforms another?
 we care only about the **sign** of ΔDCG_i

$$\Delta DCG_i = DCG_{r_1} - DCG_{r_2} = \sum_{j=1}^N g_{r_1, j} - g_{r_2, j}$$

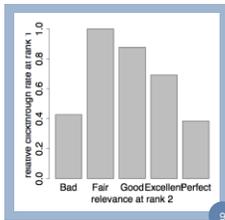
$$g_{r_j} = \begin{cases} rel_i & r_j(i) = 1 \\ \frac{rel_i}{\log_2 r_j(i)} & 1 < r_j(i) \leq l \\ 0 & r_j(i) > l \end{cases}$$

$r_j(i)$ is the rank of document i from system j

- Compute $P(\Delta DCG_i < 0)$
 - Monte Carlo simulation
 - Draw samples according to $p(X_i, X_2, \dots, X_l | q, c)$

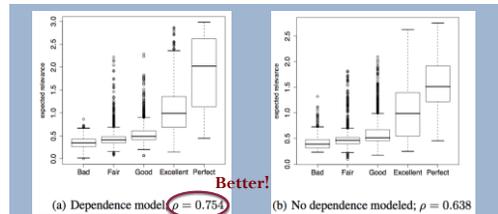
Experiments

- Data
 - From Yahoo!
 - 28,961 relevance judgments for 2021 queries
 - A distinct list includes:
 - a query
 - the ranks of retrieved advertisements
 - The clickthrough rate at each rank
- Dependency of Clicks on Entire Result List
 e.g. $rel_1 = \text{Excellent}$



Experiments 1 -- relevance prediction model

- Compare two models
 - dependence model: $p(\mathbf{X} | q, c) = \prod p(X_i | q, c)$
 - independence model: $p(\mathbf{X} | q, c) = \prod p(X_i | q, c_i)$
- Predicted relevance: $E(X_i)$



Experiments 2 -- estimating DCG

- Methodology
 - Test sets: lists with complete relevance judgments and >500 impressions (1720 lists)
 - Train sets: remaining lists with >200 impressions (>5000 lists)
 - Ground-truth: DCG
 - Predictions: $E(D\hat{C}G)$, $E(CTR) = \frac{1}{k} \sum c_k$
- Experiment 2-1
 $Cov[E(CTR), DCG] = 0.662$
 $Cov[E(D\hat{C}G), DCG] = 0.876$ **Better!**

Experiments 2 -- estimating DCG

- Experiment 2-2
 without vs. with additional two manual judgments on documents recommended by the system

Confidence in ΔDCG : $P(\Delta DCG < 0)$

Confidence	0.5 - 0.6	0.6 - 0.7	0.7 - 0.8	0.8 - 0.9	0.9 - 0.95	0.95 - 1.0
Accuracy clicks-only	0.522	0.617	0.734	0.818	-	-
Accuracy 2 judgments	0.572	0.678	0.697	0.890	0.918	0.940

Summary

- Propose a method to evaluate search engines by modeling the relationship between relevance and clicks
- Predict relevance using clicks
 - Dependence model $p(X | q, \mathbf{c}) = \prod_{i=1}^l p(X_i | q, \mathbf{c})$
- Estimate DCG with the predicted relevance
- Compare different rankings

13

Thank you !

Questions?

14

Appendix

Select Documents to Judge

- What if confident estimates are low?

Obtain more relevance judgments from human.

Intuitions:

- $r_1(i) = r_2(i)$ ignore
- $r_1(i) \gg r_2(i) \parallel r_2(i) \gg r_1(i)$ informative
- $rel_i > rel_j$

$$\Delta DCG_l = \sum_{i=1}^N g_{i1} - g_{i2}$$

$$g_{ij} = \begin{cases} rel_i & r_j(i) = 1 \\ \frac{rel_i}{\log_2 r_j(i)} & 1 < r_j(i) \leq l \\ 0 & r_j(i) > l \end{cases}$$

Algorithm 1 Iteratively select documents to judge until we have high confidence in ΔDCG .

- while** $1 - \alpha < P(\Delta DCG < 0) < \alpha$ **do**
- $i^* \leftarrow \max_i |E[G_{i1}] - E[G_{i2}]|$ for all unjudged documents i
- judge document i^*
(human annotator provides rel_{i^*})
- $P(X_{i^*} = rel_{i^*}) \leftarrow 1$
- $P(X_{i^*} \neq rel_{i^*}) \leftarrow 0$
- estimate $P(\Delta DCG)$ using Monte Carlo simulation
- end while**

15