

Offline Metrics

Fernando Diaz, Google

what are offline metrics?

system evaluation that replaces real user behavior with some proxy for that behavior.

proxies

- user models
- annotations
- log data

why use offline metrics

- faster than running a real experiment
- repeatable
- avoids exposing users to ineffective systems

example: search and recommendation

task: given an a query/context and a repository of items, find relevant information

example: search and recommendation

task: given an a **query/context** and a repository of items, find relevant information

example: search and recommendation

task: given an a query/context and a **repository of items**, find relevant information

example: search and recommendation

task: given an a query/context and a repository of items, find **relevant information**

instances

$$\mathcal{D}_q = \{\langle x, y \rangle_1, \dots, \langle x, y \rangle_n\}$$

x query-document features
(e.g. term matches,
item popularity, context)
 y relevance

all metric definition for a single query or context

metric: accuracy

$$\sum_{i=1}^n (y_i - f_{\theta}(x_i))^2$$

f document scoring function

y_i relevance label for document i

x_i query-document features for document i

example: search and recommendation

task: given an a query/context and a repository of items, find relevant information

questions

- *how is information presented?*
- *how do users consume the information?*

metric: precision

$$\text{Prec}(\mathcal{Y}^+, \mathcal{F}_\theta^+(X)) = \frac{|\mathcal{Y}^+ \cap \mathcal{F}_\theta^+(X)|}{|\mathcal{F}_\theta^+(X)|}$$

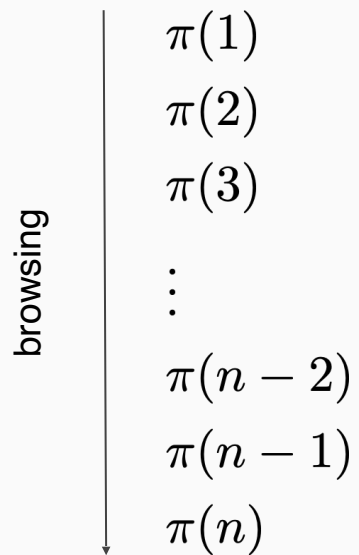
\mathcal{Y}^+ relevant document set

$\mathcal{F}_\theta^+(X)$ predicted relevant document set

metric: recall

$$\text{Rec}(\mathcal{Y}^+, \mathcal{F}_\theta^+(X)) = \frac{|\mathcal{Y}^+ \cap \mathcal{F}_\theta^+(X)|}{|\mathcal{Y}^+|}$$

ranking



metric: expected search length

user model: in-order traversal of a ranked list, collecting up to k items.

metric: number of nonrelevant documents skipped before reaching k relevant items.

$$\text{ESL}(\mathcal{Y}^+, \pi, k) = \min_{i \in \mathcal{Y}^+} \text{min-}k \pi(i)$$

min- k k th smallest value

$\pi(i)$ rank position of item i

metric: R-precision

user model: in-order traversal of a ranked list, collecting all relevant items.

metric: precision when recall is 1.

$$\text{RPrec}(\mathcal{Y}^+, \pi) = \text{Prec}(\mathcal{Y}^+, \pi_{1:k^*})$$

$$k^* = \max_{i \in \mathcal{Y}^+} \bar{\pi}(i)$$

metric: R-precision

user model: in-order traversal of a ranked list, collecting all relevant items.

metric: precision when recall is 1.

$$\begin{aligned}\text{RPrec}(\mathcal{Y}^+, \pi) &= \text{Prec}(\mathcal{Y}^+, \pi_{1:k^*}) \\ &= \frac{|\mathcal{Y}^+|}{\text{ESL}(\mathcal{Y}^+, \pi, |\mathcal{Y}^+|)}\end{aligned}$$

metric: reciprocal rank

user model: in-order traversal of a ranked list, satisfied by one item.

metric: inverse of the number of documents skipped before reaching the relevant item.

$$\text{RR}(\mathcal{Y}^+, \pi) = \max_{i \in \mathcal{Y}^+} \frac{1}{\pi(i)}$$

metric: reciprocal rank

user model: in-order traversal of a ranked list, satisfied by one item.

metric: inverse of the number of documents skipped before reaching the relevant item.

$$\begin{aligned}\text{RR}(\mathcal{Y}^+, \pi) &= \max_{i \in \mathcal{Y}^+} \frac{1}{\bar{\pi}(i)} \\ &= \frac{1}{\text{ESL}(\mathcal{Y}^+, \pi, 1)}\end{aligned}$$

metric: average precision

user model: in-order traversal of a ranked list, collecting all relevant items.

metric: precision averaged over all recall levels.

$$\text{AP}(\mathcal{Y}^+, \pi) = \frac{1}{|\mathcal{Y}^+|} \sum_{i \in \mathcal{Y}^+} \text{Prec}(\mathcal{Y}^+, \pi_{1:\bar{\pi}(i)})$$

metric: average precision

user model: in-order traversal of a ranked list, collecting all relevant items.

metric: precision averaged over all recall levels.

$$\begin{aligned} \text{AP}(\mathcal{Y}^+, \pi) &= \frac{1}{|\mathcal{Y}^+|} \sum_{i \in \mathcal{Y}^+} \text{Prec}(\mathcal{Y}^+, \pi_{1:\bar{\pi}(i)}) \\ &= \frac{1}{|\mathcal{Y}^+|} \sum_{r=1}^{|\mathcal{Y}^+|} \frac{r}{\text{ESL}(\mathcal{Y}^+, \pi, r)} \end{aligned}$$

metric: rank-biased precision

user model: in-order traversal of a ranked list; utility independent of stopping probability.

metric: expected utility given simulated user behavior.

$$\text{RBP}(y, \pi) = (1 - \gamma) \sum_{r=1}^n y_{\pi(r)} \gamma^{r-1}$$

γ patience parameter

metric: expected utility

user model: in-order traversal of a ranked list, gains utility of 1 for each relevant item.

metric: expected utility given simulated user behavior.

$$EU(y, \pi) = \sum_{r=1}^n y_{\pi(r)} \phi(y_{\pi(r)}) \underbrace{\gamma^{r-1} \prod_{r'=1}^{r-1} (1 - \phi(y_{\pi(r')}))}_{\text{examination probability}}$$

$\phi(i)$ probability that user stops given
relevance of item i

metric: time-based gain

user model: in-order traversal of a ranked list; utility independent of stopping probability; utility based on time to reach the position.

metric: expected utility given simulated user behavior.

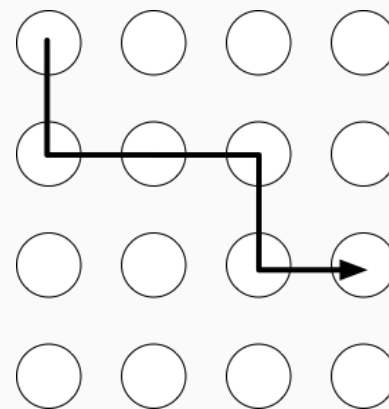
$$\text{TBG}(y, \pi) = \sum_{r=1}^n y_{\pi(r)} \Delta(y, \pi, r)$$

$\Delta(y, \pi, r)$ discount based on time to reach rank r

grid metrics

user model: model of ordered user traversal of grid items; utility independent of stopping probability.

metric: expected utility given simulated user behavior.



aggregating metrics

- traffic-weighted
- unique query/context
- usage-segmented
- group-segmented

further reading

IR evaluation: Modeling user behavior for measuring effectiveness

Charles L.A. Clarke, Mark D. Smucker, and Emine Yilmaz

SIGIR, 2015

Ian Soboroff. Building Test Collections, SIGIR 2017.

Offline evaluation options for recommender systems

Rocío Cañamares, Pablo Castells, and Alistair Moffat

Information Retrieval Journal, 2020