

Comparison of online and offline evaluation metrics in Recommender Systems

Petr Kasalický

Department of Applied Mathematics
Faculty of Information Technology
Czech Technical University in Prague

September 30, 2021

- 1 Recommender System - overview
 - Definition
 - Source data
- 2 Offline evaluation
 - Rating prediction-based evaluation
 - Ranking-based evaluation
- 3 Online evaluation
 - A/B test
 - Click-through rate

- 1 Recommender System - overview
 - Definition
 - Source data
- 2 Offline evaluation
 - Rating prediction-based evaluation
 - Ranking-based evaluation
- 3 Online evaluation
 - A/B test
 - Click-through rate

- *“The goal is to generate meaningful recommendations to users for items that might interest them” [1]*
- Requirements are more and more complex
 - Evaluation methods/metrics need to adapt

- Attributes of items and users
 - Content-based recommendation
- Interactions (feedback) between users and items
 - Explicit vs implicit
 - Interaction attributes
 - Timestamp
 - Type of interaction (click, purchase, rating)

$$F = \{f_1, \dots, f_p\}$$

$$f_j \in (U \times I \times \mathbb{Z}_t \times \mathbb{R}_v) \text{ for } \forall j \in \{1, \dots, p\}$$

- Missing-Not-A-Random (MNAR) problem
- Cold-start problem

Interaction (rating) matrix

$$F_{u,i} = \{(u_j, i_j, t_j, v_j) \mid (u_j, i_j, t_j, v_j) \in F : u_j = u \wedge i_j = i\}$$

$$r_{u,i} = \zeta(F_{u,i}), \zeta : \{0, 1\}^{|U| \times |I| \times \mathbb{Z}_t \times \mathbb{R}_v} \rightarrow \mathbb{R}?$$

	item_247	item_837	item_196	item_161	item_919	item_594	item_632		item_138			
user_5748	1			-1		0		•	•	•		1
user_3816	-0.5				0.5			•	•	•		
user_6491			-1					•	•	•		
user_8039		0.5	-0.5					•	•	•		0.25
user_2970				0.75				•	•	•		
user_6176	1					0.25		•	•	•		
user_1015			-1					•	•	•		0.5
								user's interaction vector				
			•			•				•		
			•			•				•		
			•			•				•		
user_7719		-1		0.5				•	•	•		-1

item's interaction vector

- 1 Recommender System - overview
 - Definition
 - Source data
- 2 Offline evaluation
 - Rating prediction-based evaluation
 - Ranking-based evaluation
- 3 Online evaluation
 - A/B test
 - Click-through rate

- Based on already collected data
- Fixed set of interactions split to training and test subset

$$F = F^{train} \cup F^{test} \text{ for } F^{train} \cap F^{test} = \emptyset$$

Rating prediction-based evaluation

- The task is to predict missing ratings in the interaction table

$$\tau : \mathbb{R}_?^{|U| \times |I|} \rightarrow \mathbb{R}^{|U| \times |I|}$$

- Predicted ratings:

$$T = \{(\hat{r}_{u,i}, r_{u,i}) \mid i \in I, u \in U : r_{u,i} \neq ? \wedge r_{u,i} \in R^{\text{test}}\}$$

- MSE:

$$E^{\text{MSE}}(T) = \frac{\sum_{\hat{r}_{u,i}, r_{u,i} \in T} (\hat{r}_{u,i} - r_{u,i})^2}{|T|}$$

- Similarly RMSE, MAE, NMAE etc.
- used until about 2010 when several experiments showed that minimising RSME may not improve the quality of recommendation algorithms in practical applications

Ranking-based evaluation

- Classification of relevant and not-relevant items for each user:

	classified as relevant	classified as not-relevant
truly relevant	true-relevant (TR)	false-not-relevant (FN)
truly not-relevant	false-relevant (FR)	true-not-relevant (TN)

- and measuring:

$$precision = \frac{TR}{TR + FR} \quad \text{and} \quad recall = \frac{TR}{TR + FN}$$

$$F\text{-score}_\gamma = (1 + \gamma^2) \frac{precision \times recall}{(\gamma^2 \times precision) + recall}$$

- Which items are truly relevant for a particular user?

Top- K recommendation

- Only K items can be labeled as relevant for one recommendation
- Most common scenario in practical applications
- Recall and other metrics are typically measured on subset of users:

$$U = U^{train} \cup U^{test} \text{ for } U^{train} \cap U^{test} = \emptyset$$

$$F^{train} = \bigcup_{\substack{u \in U^{train} \\ i \in I}} F_{u,i}$$

$$F^{test} = \bigcup_{\substack{u \in U^{test} \\ i \in I}} F_{u,i}$$

Recall measured using leave-one-out validation

- One item is hidden and it should be recommended based on the other items interacted by the user
- Does not follow the user's behavior over time, disadvantages *remainder* model

$$\text{recall}@K_{LOO} = \frac{\sum_{u \in U^{\text{test}}} \sum_{i \in RI_u^{\text{test}}} |\{i\} \cap \text{Top}(K, RI_u^{\text{test}} \setminus \{i\})|}{\sum_{u \in U^{\text{test}}} |RI_u^{\text{test}}|}$$

- Protection against bot and long tailedness: user-normalization:

$$\text{recall}@K_{LOO}^{UN} = \sum_{u \in U^{\text{test}}} \frac{\sum_{i \in RI_u^{\text{test}}} |\{i\} \cap \text{Top}(K, RI_u^{\text{test}} \setminus \{i\})|}{|RI_u^{\text{test}}|}$$

Recall measured using leave-last-one-out validation

- Computationally more expensive than LLOO
- Aggregation of interactions in sliding time window is recommended

$$\text{recall@}K_{LLOO} = \frac{\sum_{\substack{u \in U^{\text{test}} \\ (i_1, t_1) \in F_u^{\text{test}}}} |\{i_1\} \cap \text{Top}(K, \{i_2 \mid (i_2, t_2) \in F_u : t_2 < t_1\})|}{\sum_{u \in U^{\text{test}}} |RI_u^{\text{test}}|}$$

$$F_u = \{(i_j, t_j) \mid (u_j, i_j, t_j, v_j) \in F : u_j = u\}$$

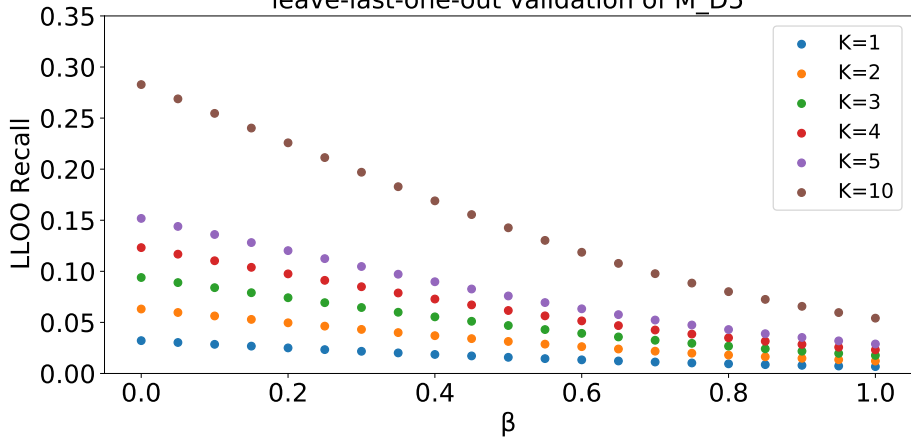
Popularity-stratified LOO recall [2]

$$\text{recall@}K_{LOO,PS}^{\beta,UN} = \sum_{u \in U^{\text{test}}} w^{\beta}(u) \frac{\sum_{i \in RI_u^{\text{test}}} |\{i\} \cap \text{Top}(K, RI_u^{\text{test}} \setminus \{i\})| p(i)^{-\beta}}{\sum_{i \in RI_u^{\text{test}}} p(i)^{-\beta}}$$

$$p(i) = \frac{\sum_{u \in U^{\text{train}}} |F_{u,i}|}{\sum_{j \in I} \sum_{u \in U^{\text{train}}} |F_{u,j}|}$$

- $\beta \in [0, 1]$ determines how much popular items should be penalised

leave-last-one-out validation of M_D5



- 1 Recommender System - overview
 - Definition
 - Source data
- 2 Offline evaluation
 - Rating prediction-based evaluation
 - Ranking-based evaluation
- 3 Online evaluation
 - A/B test
 - Click-through rate

- Implicit vs explicit
- Typically expensive to measure (fake customers, engineers effort)
- Only few models can be measured
- Unrepeatable
- Examples: CTR, Conversion rate, Customer lifetime value

- Common way to compare different versions of the system
- Users are divided into groups
- Each user group is presented with a different recommendation algorithm
- Necessary to be aware of robots, scrapers, and other nonhuman users

$$CTR = \frac{\text{number of accepted recommendations}}{\text{number of recommendations}}$$

- An example of a short-term reward
- Significantly influenced by other factors

Implicit Click-through rate (iCTR)

- Hyperparameter d - size of the time window

$$REC : Z_t \times U \times \{0, 1\}^I$$

$$iCTR(d) = \frac{\sum_{(t,u,I') \in REC} \text{sgn}(|I' \cap F_u(t, d)|)}{|REC|}$$

$$F_u(t, d) = \{i_j \mid (i_j, t_j) \in F_u : (t_j \geq t) \wedge (t + d \geq t_j)\}$$

- 1 Melville P., Sindhvani V. (2011) Recommender Systems. In: Sammut C., Webb G.I. (eds) Encyclopedia of Machine Learning. Springer, Boston, MA.
- 2 Harald Steck. 2011. Item popularity and recommendation accuracy. In Proceedings of the fifth ACM conference on Recommender systems (RecSys '11). Association for Computing Machinery, New York, NY, USA, 125–132.