May 2020

# End-to-end Optimization of Multistage Recommender Systems

Anonymous

## End-to-end Optimization of Multistage Recommender Systems
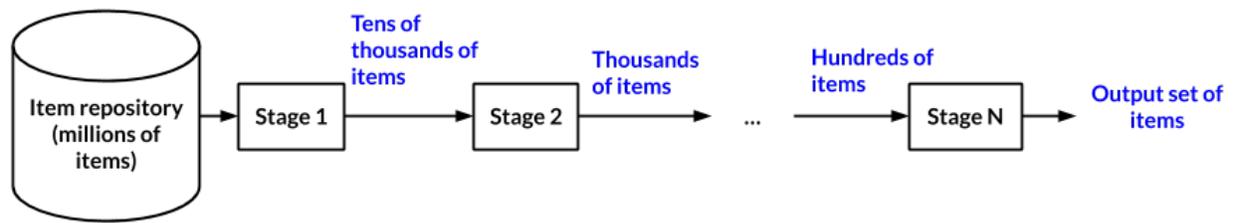
**ABSTRACT**

Complex multistage recommender systems that utilize multiple stages of ranking models and non-model parameters at different stages are used to identify a set of items as output, e.g., advertisements to be delivered by an advertising network. This disclosure describes a framework to optimize the non-model parameters to improve recall, defined based on items delivered as output in comparison to groundtruth items. The optimization can be performed offline, using a simulation that takes as input candidate items and labels of items that are known positives. The optimization can improve the quality of recommendations and can reduce computational cost.

**KEYWORDS**

- Recommender system
- Item ranking
- Non-model parameters
- Recall optimization
- Hard recall
- Soft recall
- Online advertising

**BACKGROUND**

Complex multistage recommender systems are used in various contexts to select items to recommend. For example, advertisements displayed online, e.g., on webpages, social media platforms, image sharing platforms, etc. are selected with the use of multistage recommender systems. In many contexts, the recommended items may be a small number, e.g., 1 item, 3 items, 5 items, etc. while the total universe of items to select from may include millions of items.
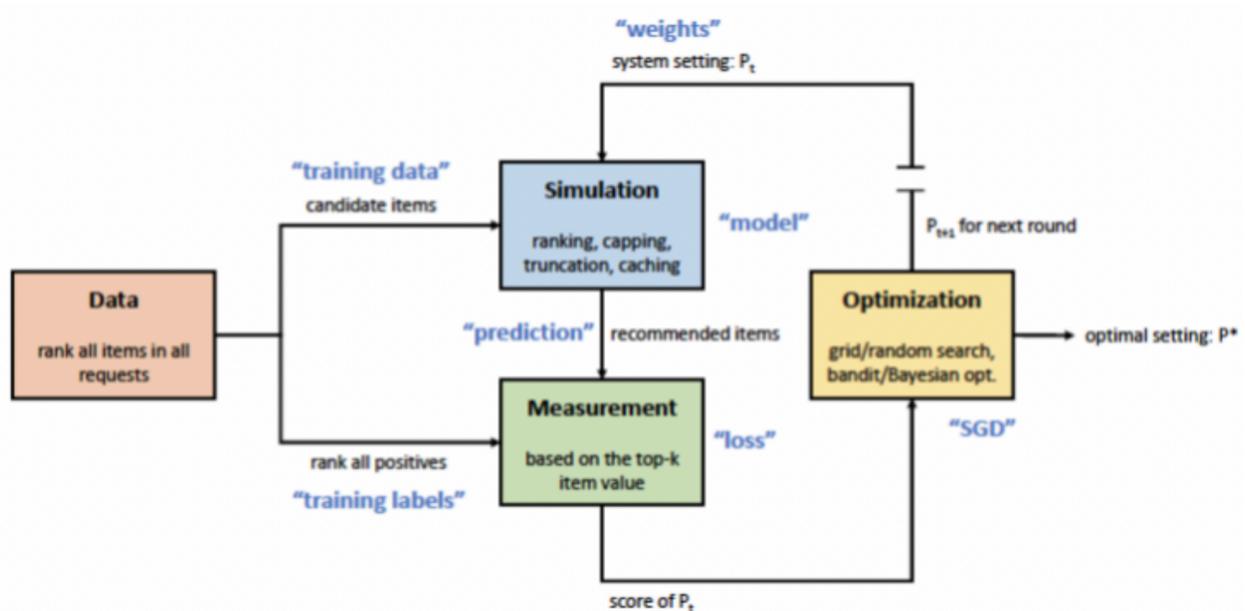
**Fig. 1: Multistage recommender system**

Multistage recommender systems use ranking models that rank items and also have additional parameters that are utilized in the selection of items to be provided. Such parameters can serve many purposes such as balancing computational resources between the various stages of the multistage recommender system, enforcing business rules (e.g., the number of items from the same provider), defining settings such as cache size and freshness used for caching of items, etc. Parameter settings can relate to various aspects of the multistage recommender system. For example, a capping process may be performed at any of the ranking stages; a truncation process may be used to limit the number of items at each stage; a caching mechanism may be used to reuse results from prior requests to the recommender system; etc. The parameters ensure that the selection process of items meets the constrained objectives defined by the value of the final set of items and the computational resources utilized to determine the final set. Such parameter settings may be selected manually (or via suboptimal mechanisms), are not easy to monitor or update, and do not guarantee the performance of the overall system in providing valuable recommendations.

One or more of the stages of the multistage recommender system may utilize ranking models that can be trained online (e.g., using A/B testing) or offline (e.g., using a loss function and stochastic gradient descent). However, optimization of non-model parameters is difficult to scale and can take significant time and resources.

## DESCRIPTION

This disclosure describes techniques for end-to-end optimization of multistage recommender systems that include ranking models and parameters. The resultant system can identify items with good recall (useful items have a high likelihood of being included in the final set) with low computational cost. For example, when the items are advertisements, the techniques can maximize value for advertisers (that pay to deliver advertisements) and the platform that delivers the advertisements.

**Fig. 2: Framework for parameter optimization of a multistage recommender system**

Fig. 2 illustrates an example framework for parameter optimization of a multistage recommender system, per techniques described herein. Data is obtained, e.g., from a late-stage ranking model, that ranks items in data in requests for recommendations received by the multistage recommender system. The multistage recommender system is run in a simulation mode with an initial set of values of the various parameters ("$P_t$") for different stages of the

system. In the example shown in Fig. 2, the parameters are related to ranking of items, capping, truncation, and caching of items at various stages of the multistage recommender system.

Candidate items are provided as input ("training data") for the simulation and a set of recommended items is obtained as the output ("prediction"). Measurement of the quality of the prediction is performed based on the top-k item value and based on ranks generated by the late-stage ranking model ("training labels"). A score is output for the parameters of the multistage recommender system and is provided as input for optimization. Optimization is performed using, e.g., stochastic gradient descent. Any suitable optimization technique such as grid search, random search, alternating optimization, Bayesian optimization, multi-arm bandit, etc. can be used. Most optimization strategies may require multiple iterations of simulation and measurement to obtain optimal parameter values.

After optimization, a new set of parameters ("$P_{t+1}$") is obtained. For example, the parameters can include weights for various settings. The new set of parameters is used for another round of simulation. Multiple such rounds can be performed until an optimal setting of parameters ("P*") is found. For example, a recall metric can be used for optimization.

For example, the output items at a particular round of simulation may be denoted by **D** and the groundtruth items (generated by the late-stage model) may be denoted by **G**. Items **D** may be generated by using candidate items **C** as input to the multistage recommender system configured with parameter settings $\mathbf{P_t}$ where **t** denotes the iteration of the simulation, such that **D=simulate(C,P)**. A recall metric **R** may then be defined as **R = measure(D, G)**. The recall measurement function can be defined in different ways. For example, a hard recall measure may be based on a ratio of the number of items that are in both **D** and **G** to the total number of items

in **G**. In another example, a soft recall measure may be defined based on a value of the items (e.g., where different items may have different values, such as in the advertising context), as the total value of items in **D** divided by the total value of items in **G**. Optimization can then be performed by finding parameter values that maximize recall - **P\*=argmax(R)**.

For example, capping parameters can be selected based on the keys associated with items in the set of items such that a limited number of items per key are retained at different stages, thus enhancing the diversity of items recommended by the system. Larger caps can be used in earlier stages of the system. Optimization can include evaluation of new capping keys, or exploration of the consistency of the capping process across multiple stages.

In another example, truncation parameters can determine the number of items that are passed from one stage to the subsequent stage of the multistage recommender system. Optimization of truncation parameters can be based on computational resources utilized with different truncation levels and scores of items that are passed.

A caching mechanism can be used at different stages to save computational resources. Caching can be based on item keys that may be associated with a time-to-live (TTL) value in the cache. Caching can be optimized based on variations in the use of the items delivered as output of the multistage recommender system, e.g., across different users. For example, the cache can be expanded to include positives.

The simulation can be performed for a given setting of non-model components, or for multiple settings in combination. The computational cost of the optimization techniques described herein is low since only the non-model components need to be simulated. This can allow for fast iteration and exploration using various parameter values. While the foregoing

description refers to use of a late-stage model to identify items that are groundtruth positives, the optimization technique can be used with other techniques to identify such items.

**CONCLUSION**

Complex multistage recommender systems that utilize multiple stages of ranking models and non-model parameters at different stages are used to identify a set of items as output, e.g., advertisements to be delivered by an advertising network. This disclosure describes a framework to optimize the non-model parameters to improve recall, defined based on items delivered as output in comparison to groundtruth items. The optimization can be performed offline, using a simulation that takes as input candidate items and labels of items that are known positives. The resultant system can identify items with good recall (useful items have a high likelihood of being included in the final set) with low computational cost. optimization can improve the quality of recommendations and can reduce computational cost.

**REFERENCES**

1. Letham, Benjamin, Brian Karrer, Guilherme Ottoni, and Eytan Bakshy. "Constrained Bayesian optimization with noisy experiments." *Bayesian Analysis* 14, no. 2 (2019): 495-519.

2. Dimmery, Drew, Eytan Bakshy, and Jasjeet Sekhon. "Shrinkage Estimators in Online Experiments." In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*, pp. 2914-2922. 2019.