

Offline Evaluation of Recommender Systems in a User Interface with Multiple Carousels

Maurizio Ferrari Dacrema^{1,*}, Nicolò Felicioni¹, Paolo Cremonesi^{1,2}

¹Department of Electronics Information and Bioengineering, Politecnico di Milano, Milano, Italy

²ContentWise, Milano, Italy

Correspondence*:

Maurizio Ferrari Dacrema

maurizio.ferrari@polimi.it

2 ABSTRACT

Many video-on-demand and music streaming services provide the user with a page consisting of several recommendation lists, i.e., *widgets* or *swipeable carousels*, each built with a specific criteria (e.g., most recent, TV series, etc.). Finding efficient strategies to select which carousels to display is an active research topic of great industrial interest. In this setting the overall quality of the recommendations of a new algorithm cannot be assessed by measuring solely its individual recommendation quality. Rather, it should be evaluated in a context where other recommendation lists are already available, to account for how they complement each other. The traditional offline evaluation protocol however does not take this into account. To address this limitation, we propose an *offline evaluation protocol for a carousel setting* in which the recommendation quality of a model is measured by how much it improves upon that of an already available set of carousels. We also propose to extend ranking metrics to the two-dimensional carousel setting in order to account for a known position bias, i.e., users will not explore the lists sequentially, but rather concentrate on the top-left corner of the screen. Finally, we describe and evaluate two strategies for the ranking of carousels in a scenario where the technique used to generate the two-dimensional layout is agnostic on the algorithms used to generate each carousel. We report experiments on publicly available datasets on the movie domain to show how the relative effectiveness of several recommendation models compares. Our results indicate that under a carousel setting the ranking of the algorithms changes sometimes significantly. Furthermore, when selecting the optimal carousel layout accounting for the two dimensional layout of the user interface leads to very different selections.

Keywords: Recommender Systems, User Interface, Evaluation, Layout Optimization, Industrial Scenario

1 INTRODUCTION

The general goal of a recommendation system is to help the users navigate the large number of options at their disposal by suggesting a limited number of relevant results. Traditionally, the focus of newly developed recommendation systems is to generate the best possible ranked list of results, see (Herlocker et al., 2004; Cremonesi et al., 2010; Sanderson and Croft, 2012). A common assumption in almost all research works is that the recommendations will be provided to the user as a single list which will be explored following its order from the first element to the last.

However, many industrial applications provide users with a two-dimensional layout of recommendations. Examples are video on demand (e.g., Netflix, Amazon Prime Video) and music streaming services (e.g., Spotify). In these scenarios the user is provided with an interface composed of multiple rows, each row containing thematically consistent recommendations, e.g., most recent, most popular, editorially curated (see Figure 1), see (Wu et al., 2016; Elahi and Chandrashekar, 2020; Pérez Maurera et al., 2020; Gruson et al., 2019; Bendada et al., 2020). These rows are referred to as *widgets*, *shelves* or as *carousels*. In a carousel interface scenario the user satisfaction depends both on the entire set of carousels shown to the user, rather than on a single list, and on their relative positions. Although finding appropriate combinations of algorithms and ranking them to provide the user with a personalized page is an active research topic of significant industrial interest (Bendada et al., 2020; Ding et al., 2019; Wu et al., 2016), the amount of work done on the design and evaluation of recommender systems based on multiple carousels is limited to a few papers based on proprietary datasets.

The two-dimensional layout is becoming more and more common thanks to *content aggregators*, i.e., services or platforms that aggregate within a single UI the content (and recommendations) provided by different services, e.g., Android TV, Sky Q, Prime Video. Although two dimensional layouts are not new, their semantics have changed over time. The first two-dimensional interfaces used a grid layout which differs from the currently used carousel layout because the rows did not have clearly identifiable differences in their semantics. The first to introduce the carousel interface was Netflix, mainly because the carousel interface allows to label the recommendation lists in a way that is much more explainable, e.g., a user that has seen the movie *Dune* may be provided with a carousel labeled "Movies similar to *Dune*". The carousel layout has several characteristics that distinguish it from the traditional single-list scenario and other grid layouts. Three important factors that should be taken into account are: (i) different carousels are generated by different and independent machine learning pipelines; (ii) some carousels are editorially curated and cannot be modified, business constraints may play a role in what contents should be promoted; (iii) content aggregators show carousels that are generated by different content providers and the content aggregator can personalize the layout but has no control on the recommendation lists content.

Among the challenges researchers face is the absence of a standardized evaluation procedure that accounts for these factors and for how the users navigate a two dimensional interface. One of the strategies that are adopted to evaluate offline in a carousel setting is to create a single recommendation list which concatenates all carousels provided to the user (Gruson et al., 2019). This strategy is however not realistic as it assumes that all recommendation lists will be centrally collected and processed (e.g., removing duplicates) and that the user explores the two-dimensional user interface sequentially, both of which are not realistic assumptions. In reality, due to strict real-time requirements or business constraints the recommendation lists cannot be collected and modified in a centralized post-processing step. Furthermore, users tend to start from the top-left corner of the screen and then proceed to explore the items both to the right and to the bottom (Kammerer and Gerjets, 2010; Zhao et al., 2016). This behavior has been long known and is referred to as "golden triangle" or "F-pattern". An example taken from an information retrieval application (Chierichetti et al., 2011) is shown in Figure 3.

In this paper we address the problem of evaluating a recommender systems in which the user interface is composed by multiple and independent carousels. This problem is referred to as the page personalization problem. This work extends prior works by (Felicioni et al., 2021b,a) and has the following contributions:

- The characteristics of a carousel interface are defined in detail as well as the relevant tasks;

- An evaluation procedure is proposed that accounts for the characteristics of the carousel interface and adapts commonly used accuracy, ranking and beyond-accuracy metrics;
- An extension of the widely used NDCG ranking metric is proposed to account for how users navigate a two dimensional user interface;
- Two simple greedy strategies are proposed as baselines for the task of carousel ranking; the strategies are agnostic on how each carousel has been created;
- The experimental analysis shows the impact of accounting for the carousel interface and two dimensional layout, yielding to sometimes very different results compared to traditional single-list evaluation approaches.

2 RELATED WORK

Most of the works targeting a carousel user interface come from industrial research with evaluation performed through A/B testing. This demonstrates the importance, at the industrial level, of identifying an adequate combination of carousels to show to the user. However, it may also indicate that the scarcity of appropriate datasets and the lack of a standardized offline evaluation protocol is hampering researchers who do not have easy access to online evaluation infrastructure, preventing them to work on the topic.

Wu et al. (2016) analyze the problem of optimizing the position of the carousels on the interface of Netflix, a popular online video streaming service. The authors propose a model based on the notion that the benefit of recommending a certain item depends on how novel it is relatively to the previous recommendations, a concept referred to as *submodularity*. Their algorithm is able to leverage scrolls and navigation feedback to dynamically optimize the user interface. The evaluation is done both online and offline. For the offline evaluation standard metrics are reported (MRR and Precision) considering a carousel as a single item which is relevant if it contains at least a relevant recommendation, therefore not accounting for the ranking within the carousel. Gruson et al. (2019) optimize the homepage of Spotify, an online music streaming service, which recommends playlists in thematic rows referred to as *shelves*. The paper evaluates a series of policies to rank the most relevant playlists for a user. The policies are ranked according to first an online evaluation and then an offline evaluation. The two rankings are then compared and de-biasing techniques are discussed to improve the correspondence of offline evaluations with online ones. In the offline evaluation, the carousels are sequentially concatenated as a single long recommendation list. Bendada et al. (2020) propose a contextual multi-armed bandit online approach to optimize the two-dimensional layout of the homepage of the online music streaming service Deezer. Each user is shown a set of carousels which the user can swipe to reveal further items. The system does not log the swipe actions. To estimate swipe actions the authors rely on the *cascade model* (Craswell et al., 2008) assuming the user has seen all items before the one they interacted with and therefore has swiped and ignored them. The policies are evaluated both online and offline with expected cumulative regrets. Ding et al. (2019) targets the problem of whole page optimization for the homepage of Amazon Video, a video streaming service. They assume that a set of widgets is already available and that the objective is to select which widgets to show and in which order while also accounting for the business constraints of the homepage. The task of carousel selection was also tackled by (Ferrari Dacrema et al., 2021) which proposed a quadratic optimization problem to decide which carousels should be selected among a pool of available ones. The optimization problem was proposed in a formulation that allowed it to be solvable using quantum computing technologies. (Felicioni et al., 2021a) proposed an extension to NDCG that accounts for the two-dimensional structure of the user interface with a different ranking discount term.

To the best of our knowledge Elahi and Chandrashekar (2020) is the only work that try to account for the two-dimensional structure of the page during the algorithm training phase. However, this scenario is unrealistic in almost all industrial settings, as different carousels are generated independently and the service that manages the layout of the page has no control on the content of each carousel.

It can be seen from these works that the carousel scenario is not yet approached in a homogeneous way and that different papers account for different factors and evaluate in different ways. Still, some aspect of the carousel scenario emerge such as the need to account for how carousels complement each other (i.e., avoid duplicate recommendations), the desire to minimize the user actions required, the tasks of carousel selection and ranking, and the lack of a common evaluation procedure.

3 CHARACTERISTICS OF A CAROUSEL INTERFACE

The carousel interface layout and the way it is usually generated by video-on-demand and music streaming platforms has important characteristics that distinguish it from a single-list setup, see (Felicioni et al., 2021a). While a carousel layout may seem similar to a traditional merge-list ensemble, where multiple recommendation lists are combined into one, this is not the case. In a real scenario multiple constraints play a role and must be taken into account: (i) the two dimensional structure of the layout; (ii) how the recommendation lists are created, (iii) how the users interact with such interfaces.

- **Layout Structure:** The two dimensional user interface of almost all devices (TV sets, personal computers, tablets) is organized with multiple horizontal carousels, where each carousel is generated according to a certain (often explainable) criteria e.g., most recent, most popular, because you watched, editorially curated (see Figure 1). Some carousels or recommendations may be hidden due to limited page size and be accessible only via user actions (i.e., vertical and horizontal swipe).
- **Recommendation Lists:** The lists shown to the users within each carousel are generated with different algorithms or by different providers and independently from each other (i.e., each algorithm or provider is not aware of the existence of the other lists or of their content). Consider for example *content aggregators*, which combine carousels from different providers, e.g., Sky, Youtube, Netflix, Prime Video, etc. No single post-processing step is applied, e.g., to remove items duplicated across different carousels. This is due to various reasons: (i) it is difficult to perform further processing under the strict real-time requirements of a recommendation system; (ii) the business constraints a content aggregator is subject to may prevent changes to the recommendation lists that are generated by the providers; (iii) some carousels are manually curated and contain a fixed recommendation list. Due to this, while each individual recommendation list does not contain duplicates, the same item may be recommended in multiple carousels. The Netflix homepage shown in Figure 1 contains the TV series *Space Force* both in the *TV Comedies* and *New Releases* carousels, which is an example of case where duplicates are not removed.
- **User Behavior:** The users will focus on the top-left triangle of the screen rather than exploring the carousels sequentially. This is usually called *golden triangle*, see Figure 3. Furthermore, they will explore the recommendations in different ways according to which actions they need to perform in order to reveal them. Usually users tend to navigate more easily with simple swipes rather than repeated mouse clicks, hence their behavior, as it is known, will change according to the device (e.g., personal computer, smartphone, tablet, Smart TV).

3.1 Tasks

Within the context of a carousel layout two orders of tasks are relevant, one is the traditional development of effective recommendation algorithms that can be used to populate a single carousel, while the second is to create the two-dimensional interface layout by deciding which carousels should be displayed and where. While developing new and more effective recommendation algorithms has been a core research topic for two decades, much less attention has been paid to the second problem. Therefore, this paper focuses in particular on how to create the carousel layout in a scenario where the technique used to generate the layout is agnostic to how the recommendation lists are created. Three sub-tasks of particular importance can be defined:

- **Carousel Selection:** given a set of recommendation lists, the task is to decide which subset of them to choose. As an example consider a scenario where we want to provide the user with a carousel related to sports, and the goal is to decide *which* sports to include based on the user preferences. Similarly, there may be recommendation lists for new releases of specific genres, e.g., comedies, fantasy, science fiction, drama etc, and again the goal is to decide which ones to show. A crucial challenge of this task is that it requires to take into account how the recommendation lists complement each other Wu et al. (2016). Clearly, the aim is to ensure the user is provided with the best possible set of recommendations, but it is not beneficial to provide to the user the same recommendations in multiple carousels. Due to this, it is not sufficient to simply select the recommendation lists that have the best accuracy when measured independently because they may be redundant. Furthermore, different users or different groups of users may have different optimal layouts. Certain constraints may apply on the types of carousels, for example one may wish that there should always be a popularity-based non personalized algorithm, an editorially curated one, a genre-specific one and a further personalized model.
- **Carousel Ranking:** given an already defined set of carousels, the task is to decide in which order to display them. A carousel ranking task is equivalent to start from a default ranking, e.g., editorially curated, and develop a re-ranking strategy. The goal is to choose the ordering that allows the user to rapidly see the most relevant recommendations. In a similar way as recommendations are ranked within an individual list, in a two dimensional layout carousels should be ranked according to both (i) the user navigation pattern as it emerges from the golden triangle and (ii) the user actions required to access certain positions. The challenge of this task is to account for how the user navigates a two dimensional user interface as well as what portion of the recommendation lists are immediately visible to the user and which require user actions. Due to this, ordering the carousels according to their decreasing recommendation quality will lead to suboptimal layouts. Furthermore, business constraints may apply on the relative ordering of some of the carousels.
- **Carousel Insertion:** given an already available carousel interface, the task is to decide where a new carousel should be added without altering the relative position of the available ones. As opposed to carousel selection and ranking, the insertion of a carousel is an incremental task and does not require to search for an optimal layout anew. As an example, consider a content aggregator that displays recommendations provided by several video on demand services the user has subscribed to. If the user subscribes to a new service, we want to include that new carousel in the already existing layout, in a position that maximizes the user satisfaction within possible business constraints.

All the described tasks assume that the recommendation lists that can be used as carousels are already available and the goal is to decide what is the optimal layout. This requires to develop an evaluation

194 procedure that is tailored to the specific nature of the carousel setting as well as extend traditional
195 evaluation metrics, see Section 4.

196 As a last note, while this paper considers the creation of the carousel layout given that (i) all carousels are
197 already available and (ii) the layout manager is not aware of the algorithms used to fill the carousels, it is
198 indeed possible to combine the carousel ranking tasks at the page level with the item ranking task within
199 each carousel.

4 EVALUATION PROCEDURE FOR A CAROUSEL INTERFACE

200 This section describes an evaluation procedure tailored for a carousel user interface, which accounts
201 for its characteristics and the constraints that are normally present in a real industrial setting. While the
202 traditional evaluation assess the recommendation quality of a single recommendation model, in a carousel
203 scenario the goal of the evaluation is to assess the recommendation quality of a certain layout composed of
204 recommendation lists (either fixed or generated by a specific recommendation model). The recommendation
205 quality of a layout can then be used to compare different strategies to generate it, e.g., to rank carousels or
206 to choose which recommendation model should generate a specific carousel. A summary of the notation
207 used is reported in Table 1.

- 208 • **Scenario:** The carousel evaluation setting is closely tied to a specific user interface commonly adopted
209 for video on demand and music streaming services. As such, while it could be applied for any domain,
210 we recommend its use be limited to those that employ such carousel interface. Researchers have ample
211 freedom to decide how to generate the carousels (editorially curated lists, recommendation models,
212 etc.).
- 213 • **Optimization:** If some of the carousels are generated with recommendation models, the first step is
214 to ensure that all models are adequately optimized. In a traditional offline evaluation each model is
215 optimized independently by selecting the hyperparameters that optimize its recommendation quality
216 on a validation set (Ferrari Dacrema et al., 2019). The same holds for a carousel scenario. In a carousel
217 setting the recommendation lists provided to the user may come from different sources, e.g., a third
218 party providing the recommendation engine for some carousels, while the platform provides others,
219 (Pérez Maurera et al., 2020), the carousels shown to the user may change dynamically between sessions
220 and different users may see a different carousel selection. This means that during the optimization of
221 the models the layout of the user interface is, in general, not known and cannot be used during the
222 optimization phase. For this reason, we recommend all models should be optimized independently
223 and, if assumptions are made on the composition of the carousels, those should be stated clearly. For
224 example, it may be assumed that a popularity-based carousel is present and train recommendation
225 models specifically to improve the recommendation quality on less popular items.
- 226 • **Recommendation:** The recommendations that will be shown to the users are the sequence of all the
227 recommendation lists in the layout. Which lists to include and in which order is part of the experimental
228 setting and will be chosen according to the scenario of interest. All lists have the same length, H ,
229 corresponding to the horizontal dimension of the interface (i.e., number of columns). Given V the
230 number of recommendation lists that are displayed (i.e., the number of rows) the user will receive
231 $l = HV$ recommendations. Usually the recommendations come from different providers and there will
232 be no centralized postprocessing done on the set of all recommendation lists. While a carousel will
233 only contain a certain item at most once, there could be items appearing in more than one carousel.

- **Evaluation:** In a carousel setting the recommendations provided to the user will be displayed with a two-dimensional pattern. A frequent simplification is to concatenate all V recommendation lists in a single one of length $l = HV$ and remove duplicate recommendations. While this allows to use traditional metrics (e.g., NDCG, MAP), it makes assumptions that are not consistent with a carousel layout: (i) the user explores the lists sequentially; (ii) the recommendation lists are centrally collected and postprocessed. In reality, the user will not explore a list sequentially, but rather start from the top left corner and will move in both directions exploring multiple carousels (see Figure 2). There may also be items that appear in more than one carousel. When this happens, we must ensure that a correct recommendation is only counted once and with the correct ranking, despite having appeared multiple times in different positions. The correct recommendation should be counted where the user would *see it first*, according to the user's navigation pattern. A detailed description of how the evaluation metrics should be computed is provided in Section 4.1.

4.1 Evaluation metrics

Evaluation in a carousel setting presents broad similarities with a traditional top-n recommendation scenario. An important difference is the presence of duplicates in the recommendation list and the two-dimensional way users explore the interface, which impacts how ranking metrics may consider the item position. This section describes how accuracy, ranking and beyond accuracy metrics should be adapted to a carousel evaluation.

4.1.1 Accuracy metrics

Accuracy metrics that are based only on the total number of correct recommendations, e.g., Precision, Recall, Hit Rate, etc. can be computed in the same way as in single-list evaluation, provided that each correct recommendation is only counted once even if it appears in multiple carousels.

4.1.2 Ranking metrics

Adapting ranking metrics to a carousel evaluation requires to account for the two-dimensional user exploration of the interface at two stages: (i) how the *ranking discount* is computed; (ii) how correct recommendations appearing in more than one carousel are considered. For traditional top-n ranking with a single recommendation list (or a merge-list embedding), duplicates are removed from the list; therefore, we assume that each recommendation list does not contain duplicates. However, in a carousel evaluation duplicates might occur between carousels. Such duplicates must not be removed in order to maintain the original positions of the recommendations and mimic the real behavior of carousel-based user interfaces. In a carousel evaluation scenario, an item in the recommendation list is relevant, i.e., a correct recommendation, if it meets two conditions:

- The item appears in that user's ground truth;
- The item has been recommended only once or, if it has duplicates, it is the one corresponding to the *best ranking*, i.e., the position where the user is likely to see it first. The question therefore becomes defining a new ranking discount that accounts for the user behavior in a carousel interface. Section 5 discusses different discount functions that are applied to NDCG but can be applied seamlessly to other metrics such as MRR.

272 4.1.3 Beyond-accuracy metrics

273 Another important dimension to be taken into account when evaluating recommendation models is their
274 diversity. There are three ways to measure the diversity of recommendations: (i) *individual* or within-list,
275 measuring how diverse are the recommendations received by a single user, e.g., based on the item features or
276 on the item popularity in the training data; (ii) *list-based*, that consider the diversity of the recommendations
277 received by couples of users; (iii) *aggregate* or distributional, that measure how unbalanced is the overall
278 distribution of recommended items at a system level.

279 In a carousel evaluation setting most beyond-accuracy metrics can be used seamlessly and computed
280 on the recommendations of all carousels concatenated in a single list and counting the occurrences of
281 duplicate items as well. This is because repetitions of the same items indicate the recommendations have
282 lower diversity. It is possible to calculate correctly several commonly used within-list beyond-accuracy
283 metrics, such as Average Popularity (the average popularity of the recommended items) and Novelty (mean
284 self-information (Zhou et al., 2010), a function of the item popularity in the training data). The same applies
285 for distributional diversity such as Item Coverage (the quota of items that have been recommended at least
286 once) and other metrics that are computed on the global number of times each item has been recommended,
287 the more unbalanced the distribution the less diverse the recommendations are: Shannon diversity, Gini
288 diversity and Herfindahl Index (Adomavicius and Kwon, 2012; Ferrari Dacrema, 2021).

289 As opposed to the previous two categories, list-based diversity metrics such as Mean inter-list diversity
290 (MIL) (Zhou et al., 2010) and the equivalent Hamming diversity behave differently. MIL measures
291 the average quota of the recommendation lists any couple of users has in common. Computing MIL
292 with his approach is however very computationally expensive as it is equivalent to compute a user-
293 based KNN similarity model. The original formulation of MIL is developed for the traditional single-
294 list recommendation setting, in which all users receive the same number of recommendations and the
295 recommendations do not contain duplicates. In this scenario, a recent proof by Ferrari Dacrema (2021)
296 shows that MIL is equivalent to the Herfindahl Index and therefore is a distributional diversity metric,
297 this means it can be computed in negligible time based on the global number of times each item has been
298 recommended. In a carousel setting, due to the presence of duplicates, the quota of recommendations two
299 users have in common becomes asymmetric. As an example, consider two users *A* and *B* that received 10
300 recommendations. If they have one item in common this means that 10% of the recommendations are equal.
301 If the common item appears twice in the recommendations received by user *A*, the similarity becomes
302 asymmetric with *A* having 20% of the recommendation list in common with *B*, while *B* having 10% of
303 the list in common with *A*. Removing the duplicates would produce erroneous results since MIL is defined
304 as an arithmetic mean and requires recommendation lists of the same length. Due to this, the equivalence
305 shown in (Ferrari Dacrema, 2021) does not hold in a carousel setting and it is an open research question
306 whether the effect of the compositions of the specific user recommendation lists would cause MIL to still
307 behave as a distributional diversity metric or not.

308 In a carousel scenario it would also be possible to define additional beyond-accuracy metrics that measure
309 the diversity *between carousels*, e.g., measure the overlap between recommendations contained in different
310 carousels shown to a single user. Defining new metrics for this scenario is however an open research
311 question that goes beyond the scope of this paper.

5 EXTENDING ONE-DIMENSIONAL NDCG

The *Discounted Cumulated Gain* (DCG), as well as its *Normalized* version (NDCG) (Järvelin and Kekäläinen, 2000, 2002), are among the most used metrics for the evaluation of ranked lists. This metrics come from the information retrieval domain and are widely used to evaluate recommendation systems. The DCG metric relies on two assumptions:

1. highly relevant results are more valuable for a user;
2. within a list of results, it is preferable to have relevant results in the first positions.

Let l be the recommendation list length, i.e., cutoff, u a user within the set of existing users U and j the position in the recommendation list (note that j does not refer to a specific item but rather to a position in the list). The DCG for user u is defined as the following discounted sum of gains:

$$DCG_u = \sum_{j=1}^l g_{uj} d_j \quad (1)$$

The gain function g_{uj} is responsible for rewarding highly relevant items in position j for user u , while the discount function d_j introduces a penalization that should increase as the position j moves towards the end of the list.

Given r_{uj} as the relevance (e.g., rating) of the item in position j of the recommendation list for user u , the most used formulations for gain and discount are $g_{uj} = 2^{r_{uj}} - 1$ and $d_j = \frac{1}{\log_2(j+1)}$. Hence DCG is computed as (Burgess et al., 2005):

$$DCG_u = \sum_{j=1}^l \frac{2^{r_{uj}} - 1}{\log_2(j+1)} \quad (2)$$

The DCG for each user is normalized by computing the ideal DCG for that same user, denoted as $IDCG_u$. While the DCG considers all items in the recommendation list, the $IDCG$ is computed assigning to each item its true relevance (i.e., the one in the test data) and therefore obtaining the best possible ranking given that user's ground truth within the available recommendation list of length l . The NDCG for user u is computed as:

$$NDCG_u = \frac{DCG_u}{IDCG_u} \quad (3)$$

Finally, the global NDCG is computed as the average of the NDCG of each user:

$$NDCG = \frac{1}{|U|} \sum_{u \in U} NDCG_u \quad (4)$$

Notice that this formulation is only one of many possible formulations for the DCG. Several other ways of rewarding and discounting results have been proposed in previous research (Kanoulas and Aslam, 2009; Zhou et al., 2014). In the following, we will start from this widely used formulation and extend it. Other types of gain and discount functions can be extended in an analogous way. We leave the analysis of different gains and discounts as future work.

Extending the DCG definition to a carousel setting requires to take into account both the way a user explores a two dimensional interface, following the golden triangle, and that only a portion of the recommendations will be immediately visible to the user and further portions will become visible following user actions, e.g., click, swipe. Furthermore, the shape and size of the user interface as well as the user behavior will change according to the device.

In a two-dimensional scenario, the standard DCG definition can be naively adapted in the following way. Let H be the horizontal dimension of the interface (i.e., the length of each carousel) and V the vertical dimension of the interface (i.e., the number of carousels). The carousels can be concatenated in a single list of length $l = V \cdot H$ items on which the standard DCG formulation can be applied. This strategy assumes that the users will explore all carousels sequentially, from the first to the last, which, as previously discussed, is not consistent to the user behavior and does not account for the interface navigation constraints. Therefore, we suggest researchers *do not* apply this strategy as it does not represent a realistic scenario.

Thus, inspired by (Järvelin and Kekäläinen, 2000), we make the following assumptions the two-dimensional DCG should meet:

1. highly relevant results are more valuable for a user;
2. a relevant result is valuable to the user only when it is first seen;
3. within a grid of results, it is preferable to have relevant results close to the top-left corner;
4. it is preferable that relevant items are immediately visible to the user or can be made visible with the least effort, e.g., the least number of user actions, given that actions of different type may require different efforts.

In order to account for this set of assumptions, we propose to extend the DCG metric as a two-dimensional DCG, $2DCG$, in the following way:

$$2DCG_u = \sum_{j=1}^V \sum_{k=1}^H g_{ujk} d_{jk} \quad (5)$$

Where $j \in [1, V]$ and $k \in [1, H]$ represent the position within the two-dimensional carousel interface. As in the one-dimensional version, the gain function is responsible for rewarding highly relevant results, according to assumptions (1) and (2). The discount function, instead, should account for the penalty related to the position and number of user actions, according to assumptions (3) and (4).

Inspired by the one-dimensional version, we define $g_{ujk} = 2^{r_{ujk}} - 1$ where r_{ujk} is the relevance of item in carousel j , position k , for user u . The discount term will depend on the position in the layout, allowing ample freedom on how to define it in different use cases.

The normalized version of this metric, $N2DCG_u$ requires to define the $I2DCG_u$, which is the the $2DCG_u$ of the *ideal ranking*. In a single list setting the ideal ranking is the list which contains the relevant items in decreasing relevance from the beginning of the list. In the generalized two-dimensional layout it contains the user's most relevant items, ranked according to decreasing relevance in positions with decreasing position discount. Note that, as done for IDCG, if a user has a number of relevant items that exceed the recommendation list length l , only the l most relevant ones will be used to compute the $I2DCG_u$. The ideal ranking meets the following constraints: $g_{ujk} \geq g_{uxy}$ if $d_{jk} < d_{zx}$ for any pair of items in positions (j, k) and (x, y) .

Finally, $N2DCG$ is computed as the average of all $N2DCG_u$ as in the traditional NDCG.

5.1 Relevance

As stated in assumption (2), a relevant item is valuable for the user only when it is first encountered. This means that if a relevant item appears multiple times, each in a different carousel, it should be considered as relevant only in its *best* position. We define such position as the one with the lowest *discount*. Function $r_{u,j,k}$ should be modified accordingly. Due to this, if an item is correctly recommended twice the position that will be considered relevant may change depending on the discount function applied. An example of this is shown in Figure 5.

5.2 Single List Discount

It is possible to represent with N2DCG the scenario where all carousels are concatenated into a single list by calculating the position of cell in coordinates j, k if all carousels lists would be concatenated:

$$d_{jk}^s = \frac{1}{\log_2((j-1)H + k + 1)} \quad (6)$$

As previously mentioned, this formulation is not grounded in a realistic scenario because it does not reflect the user behavior (see Figure 5 A), therefore we argue it should not be applied.

5.3 Golden Triangle Discount

In order to account for the *golden triangle* behavior, as per assumption (3), the position discount should decrease as the distance of the cell from the top-right corner increases:

$$d_{jk}^t = \frac{1}{\log_2(\alpha j + \beta k)} \quad (7)$$

The coefficients α, β are two weights that can be used to account for different types of user behaviors. For instance, let us assume a scenario where users are more inclined to explore the vertical dimension. In this case, α should be set to a low value in order to penalize less the vertical dimension. In order to make the discount start from 1, α and β should be ≥ 1 since the base of the logarithm used is 2. Notice that this is true only because we are extending a logarithmic discount function. For other discount functions (Kanoulas and Aslam, 2009; Zhou et al., 2014) the constraints can change.

5.4 User Actions Discount

In order to account for assumption (4) the position discount should decrease the more actions are required by the user to make that position visible. In a carousel interface there is an initial rectangular portion of the recommendations that are immediately shown to the user. We refer to the number of items visible as V_h (number of visible horizontal items) and to the number of carousels visible as V_v (number of visible vertical items), see Figure 4. In order to reveal more items in either the horizontal or vertical directions, the user needs to perform a certain action, i.e., scroll the mouse wheel, swipe on mobile devices). We call this actions as *swipe actions*. Each of these actions will reveal a certain number of new items available in the recommendation lists. The number of items revealed with a swipe action is defined as $\delta_h \in \{1, 2, \dots, V_h\}$ (for horizontal swipe actions) and $\delta_v \in \{1, 2, \dots, V_v\}$ (for vertical swipe actions).

For example, on Netflix every horizontal swipe on a carousel replaces all the items displayed in the carousel, therefore $\delta_h = V_h$.

The number of horizontal and vertical actions needed to visualize an item in position (j, k) (the k -th item in carousel j) is computed as

$$n_j^{\text{h-swipes}} = \begin{cases} \left\lceil \frac{j-V_h}{\delta_h} \right\rceil & \text{if } j > V_h \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

$$n_k^{\text{v-swipes}} = \begin{cases} \left\lceil \frac{k-V_v}{\delta_v} \right\rceil & \text{if } k > V_v \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

The final discount accounts for both the triangle discount, as defined in (7), and the number of horizontal and vertical user actions required to discover the item

$$d_{jk}^a = \frac{1}{\log_2(\alpha j + \beta k + \gamma n_j^{\text{h-swipes}} + \lambda n_k^{\text{v-swipes}})} \quad (10)$$

Notice that this formulation accounts for both vertical and horizontal swipes. The coefficients $\alpha, \beta, \gamma, \lambda$ are four positive weights that can be used to account for different types of user behaviors. The first two weights (α and β) control the general penalization of the vertical and horizontal dimensions, respectively and in order for the total discount to start from 1 $\alpha, \beta \geq 1$. Penalizing more or less the user actions needed to reveal a certain item is possible by controlling $\gamma, \lambda \geq 0$. For example, it could be that items presented together in the same carousel have a similar probability of interaction (see the first 10 elements of the first carousel in Figure 2). Hence, the horizontal dimension should be penalized less. Another possibility is that, on a desktop device, the horizontal swipe done with a mouse click will have a higher weight than the same swipe done with a touch on a mobile device.

For illustrative purposes, let us consider a possible scenario for a mobile device, where the screen contains 4 carousels and 3 recommendations each. We set the horizontal and vertical steps to 1, $\alpha, \beta, \gamma, \lambda$ are set to 1 as well. The resulting discount is shown in Figure 5 (B).

Figure 6 compares the single list DCG and 2DCG in two different scenarios. The user interface has $V = 3$ carousels each of length $H = 6$, with $V = V_v = 3$, $V_h = 3$, $\delta_h = 3$ and the penalty for the horizontal user action is $\gamma = 10$. The first scenario compares the values when three correct recommendations exist, see Figure 6 (A), to which a fourth correct recommendation is added in position (3, 1), see Figure 6 (B). The value of DCG increases by 0.263 (from 1.057 to 1.320), while 2DCG increases by 0.497 (from 1.364 to 1.861) indicating how 2DCG is much more affected by the additional correct recommendation than DCG. This is because DCG applies a much higher discount due to the position compared to 2DCG. The effect will become more marked for longer carousel interfaces, making much more difficult for DCG to discriminate reliably between layouts that differ in their last recommendation lists. For a correct recommendation in position (10, 1) DCG applies a discount 40% higher than that of 2DCG, while only 20% higher if a correct recommendation is in position (30, 1). This means that although the most marked differences are in the first carousels, despite the logarithmic nature of the discount, 2DCG is still able to better account for the user behavior even in long carousel interfaces. Consider that at the time of writing the Netflix homepage shows more than 30 carousels. The second scenario considers the impact of the user action penalty on determining

the optimal layout, when three recommendation lists are available. The recommendation lists have zero, one and two correct recommendations respectively. In the list with two correct recommendations the second one is outside the visible area and requires a user action to be seen. The layout depicted in Figure 6 (C) puts first the list with two correct recommendations, then the list with one and lastly the list with zero. This corresponds to a DCG of 1.256 and 2DCG of 1.246. A different layout, obtained by swapping the first two lists as depicted in Figure 6 (D), has instead a DCG of 1.221 and 2DCG of 1.312. According to the single list DCG it is preferable to put the recommendation list with two correct recommendations at the top, while the 2DCG will account for the fact that the second correct recommendation will be outside the user visible area and will have a much lower contribution. Due to this, DCG and 2DCG will lead to the selection of different layouts. This effect will be more or less marked according to the impact of the user action penalty, which will depend on the scenario of interest.

6 EXPERIMENTAL PROTOCOL

This section describes the experimental protocol that is used for the analysis reported in Section 7.

6.1 Recommendation Lists

In a realistic carousel scenario, several recommendation lists (or carousels) are available generated with different algorithms or editorial rules. In order to mimic this setting we include in the evaluation 16 algorithms developed in the last three decades of research that are both simple, well-known and competitive (Ferrari Dacrema et al., 2021). The algorithms can be grouped into 6 categories:

- **Non-Personalized:** TopPopular, which recommends the most popular items, and GlobalEffects which recommends the highest rated items.
- **Item-Based Heuristic:** ItemKNN and UserKNN (Sarwar et al., 2001) compute the item-item or user-user similarity based on past user interactions. Both models use cosine similarity with shrinkage to reduce the similarity of items or users with low support.
- **Item-Based Machine Learning:** SLIM (Ning and Karypis, 2011), SLIM BPR and EASE^R (Steck, 2019) all compute an item-item similarity by optimizing different criteria. EASE^R in particular is a recently proposed shallow autoencoder with a closed form solution.
- **Graph-based:** P³ α (Cooper et al., 2014) and RP³ β (Paudel et al., 2017) are both simple methods that represent the user interactions as a bipartite user-item graph and model the recommendations by simulating a random walk.
- **Matrix Factorization:** PureSVD (Cremonesi et al., 2010), FunkSVD (Ferrari Dacrema et al., 2021), Non-negative matrix factorization (Cichocki and Phan, 2009), MF BPR (Rendle et al., 2009) and IALS (Hu et al., 2008) are all widely known models that represent the user interactions as the product of lower dimensionality matrices learn by optimizing different losses.
- **Content-based and Hybrid:** ItemKNN CBF builds an item-item similarity using the item features, while ItemKNN CF-CBF uses both item ratings and item features concatenated in a single vector (Mobasher et al., 2003). Both use cosine with shrinkage as the collaborative ItemKNN.

6.2 Hyperparameter optimization

While this paper does not aim to show that any particular model or carousel ranking strategy is superior to others, we nonetheless ensure that all algorithms are fairly and consistently optimized. To do so we followed the best practices highlighted by Ferrari Dacrema et al. (2021) and we relied on the framework

published. The data is split in 80% training, 10% validation and 10% test with a user-wise random holdout. (Antenucci et al., 2018; Ferrari Dacrema et al., 2021)¹ optimizing the recommendation quality on the validation split, measured with NDCG at cutoff 10. Note that in this case each model is optimized in a traditional scenario where only a single recommendation list is provided, therefore there is no need (and indeed at this stage is not possible) to account for the whole user interface. We explored 50 cases for each algorithm, the first 16 used as random initialization. We constrain the optimization to a maximum of 14 days for each recommender model and a maximum of a 64GB of RAM. Once the optimization terminates, the final model is fitted on the union of the training and validation data using the best hyperparameters found and is evaluated on the test data. The hyperparameter ranges and distribution used during the optimization are the same used and described in (Ferrari Dacrema et al., 2021). The optimal hyperparameters for all recommendation models are reported as part of the Appendix in Table 8, 9 and 10.

6.3 Datasets

We report the results for some widely known and used datasets, all of which are publicly available. We only selected datasets collected from domains that typically adopt a carousel-based user interface, i.e., video-on-demand and music streaming:

- *MovieLens 20M*² (Harper and Konstan, 2016), a popular dataset of movies recommendations. The dataset contains user provided tags for items as well as the year of release and the genre. User ratings are available in a range 1-5. The dataset has 20.0 M interactions, 27 k items, 138 k users and a density of $5.3 \cdot 10^{-3}$.
- *Netflix Prize*³ (Bennett et al., 2007), is the well known movie dataset from the *Netflix Prize*. User ratings are available in a range 1-5. The dataset has 100.4 M interactions, 17 k items, 479 k users and a density of $1.2 \cdot 10^{-2}$.
- *ContentWise Impressions* (Pérez Maurera et al., 2020) a dataset collected from a video-on-demand media provider that contains both user interactions and impressions. The interactions are implicit. The dataset has 4.5 M interactions, 145 k items, 41 k users and a density of $7.4 \cdot 10^{-4}$.

6.4 Carousel Layouts Heuristics

While the purpose of this paper is not to propose a specific algorithm to find an optimal carousel ranking, we believe it is useful to describe two simple heuristic strategies that can be used as baselines to compare with more sophisticated strategies. As previously observed, there are only few works that deal with the problem of recommendation list selection for a carousel interface. Usually, those works make assumptions on the recommendation models, leverage specific types of data such as session and context or require an online setting. In other cases, the selection of carousels is part of the recommendation model itself, which means it is not applicable in the scenario where the recommendation models are black boxes. Note that one may choose a single global layout, but also select different layouts for different users or groups of users.

Let's first consider the solution space for the task of carousel selection and ranking:

- **Exhaustive selection with default ranking:** Given a set of available recommendation lists, this baseline evaluates all possible subsets of a desired length, i.e., the number of carousels, with a default ordering that may be selected heuristically. For example, we may wish to select one among many

¹ The framework implementation relies on the Scikit-optimize package <https://scikit-optimize.github.io/>

² <https://grouplens.org/datasets/movielens/20m/>

³ <https://www.kaggle.com/netflix-inc/netflix-prize-data>

possible carousels related to sports, then another related to a genre and finally another personalized list. An exhaustive search is very computationally expensive as it corresponds to selecting the combinations without repetitions of V carousels within a pool of M lists, resulting in a total number of cases $\frac{M!}{V!(M-V)!}$. This makes an exhaustive search impossible for all but the smallest cases, where the number of available lists is very limited, hence it is of little practical use.

- **Exhaustive selection and ranking:** Given a set of available recommendation lists, the goal is not only to decide which to select, as in the *Exhaustive selection with default ranking* case, but also to decide how to rank them as carousels. This search corresponds to explore all possible permutations of the selected V carousels, which requires to evaluate $V!$ layouts for each selection, making it again too computationally expensive for practical use.

Two simple greedy strategies that can be used for the selection and ranking of carousels are the following:

- **Individual Greedy:** The recommendation lists are selected according to their recommendation quality measured individually on the validation data, and ranked with decreasing values. Therefore the list with the best recommendation quality will become the first carousel, the second-best will become the second carousel and so on until the desired number of carousels is reached. In our experiments we have used NDCG with a cutoff at 10 to rank carousels with the Individual Greedy strategy. This approach cannot account for duplicate recommendations and may select a set of lists with similar recommendations. For example, the second-best list may be very similar to the first one, so most of its correct recommendations will already be present in the first carousel. This strategy requires only the recommendation quality on the validation data which is either already available after the recommendation models that generated the lists have been optimized, or can be easily computed once a new editorially curated list is added to the available ones. Hence, the Individual Greedy selection baseline has a low computational cost.
- **Incremental Greedy:** This baseline does not select the recommendation lists based on a fixed accuracy value, but rather iteratively evaluates all of them accounting for those that have already been selected. If a recommendation list has high accuracy but provides recommendations similar to those already selected, it will exhibit lower recommendation quality. This carousel selection baseline is better suited to account for the characteristics of a carousel interface, but it is much more computationally expensive requiring to run $\sum_{i=1}^V M - i + 1$ evaluations.

7 RESULTS

In this section we apply the proposed carousel evaluation on numerous widely used algorithms and compare the results obtained with the traditional evaluation which considers each model independently, we also compare the results obtained by optimizing a carousel layout with the traditional single-list NDCG and with the proposed N2DCG. We discuss the results of this comparison and highlight some common trends.

The total number of available recommendation models, i.e., M , differs according to the dataset: Movielens 20M has 16 because it includes 14 collaborative and 2 content based ones, Netflix Prize has 14 which corresponds to all the collaborative models, ContentWise Impressions has 12 because SLIM BPR and EASE^R required more than the available 64GB of RAM and could not be optimized. N2DCG uses the following parameters: $V_v = \min(V, 3)$, $V_h = 3$, $\delta_v = 1$, $\delta_h = 3$, $\alpha = \beta = \gamma = \lambda = 1$, see Figure 5.

7.1 Recommendation quality under a carousel evaluation

The first analysis shows how the relative recommendation quality of a model changes by using a carousel evaluation. In this experiment, the goal is to choose which model to add as the *last* carousel in an interface that contains an increasing number of carousels: TopPopular, ItemKNN CF and, for the Movielens 20M dataset, ItemKNN CBF.

The models are first evaluated individually with the traditional evaluation protocol and then with the proposed carousel evaluation protocol. All recommendation lists have a length, i.e., H , of 10. Note that in the individual evaluation there will be a single recommendation list while in the carousel evaluation there will be more than one, therefore the absolute values of the NDCG are not comparable, for this reason the analysis will focus on the ranking of the models. The results are reported both as tables and as figures that highlight the changes in the relative ranking of the modes: Table 2 and Figure 7 (A) (Movielens 20M), Table 4 and Figure 7 (B) (Netflix Prize), Table 3 and Figure 7 (C) (ContentWise Impressions).

By looking at Table 2 and Figure 7 (A) (Movielens 20M) we can make several considerations. Under the traditional evaluation in which an algorithm is evaluated independently (i.e., individual) we can see that almost all personalized algorithms have a recommendation quality between two and three times better than TopPopular. The only model that is worse than TopPopular is GlobalEffects and the best performing algorithms are SLIM, UserKNN and EASE^R. If we look at the carousel evaluation, we can see that no algorithm has a lower recommendation quality than TopPopular, which is expected since TopPopular is the predefined first carousel while each algorithm is the second one and therefore the final recommendation quality can only be greater or equal. As a general trend we can see that the relative effectiveness of the models differ, resulting in changes to the ranking of the algorithms in the two evaluation modes. Some models such as GlobalEffects and PureSVD are always ranked in the same position. Others, in this case all other matrix factorization algorithms, gain 2 or 3 positions. On the other hand item-based machine learning models tend to consistently lose some positions, with EASE^R being the worst affected and losing 4 positions. As a result, while in the individual evaluation the best algorithms are SLIM ElasticNet, UserKNN CF, SLIM BPR and EASE^R, in the carousel evaluation the best algorithms are UserKNN CF, SLIM ElasticNet, IALS and FunkSVD. Since the recommendation lists generated by all algorithms are identical for both evaluation procedures, the difference in the ranking lies in how those recommendations intersect. Algorithms which will tend to recommend popular items will be penalized in this carousel evaluation because popular items will already be present in the TopPopular carousel, whereas algorithms providing accurate recommendation but involving less popular items will be advantaged. Similarly, since the second carousel layout adds an ItemKNN CF the algorithms providing similar recommendations to it will be penalized too.

Table 4 and Figure 7 (B) (Netflix Prize) shows a similar, albeit more marked, behavior. The models ranking is mostly unchanged when the recommendation model is evaluated in a carousel layout having only TopPopular as first carousels, the only changes are RP³ β and FunkSVD swapping relative positions. The inclusion of the ItemKNN CF instead causes some sharp changes in ranking, in particular EASE^R falls by 6 positions while FunkSVD gains 7. Finally, Table 3 and Figure 7 (C) (ContentWise Impressions) show a case where the ranking is less affected. Using a TopPopular as first carousel does not change at all the relative ranking of the models. Including an ItemKNN CF as second carousel affects mostly two models, P³ α which loses 4 positions and IALS that gains 3.

Overall, no clear pattern emerges, with the different recommendation models being affected in different ways according to the carousel structure and the dataset. An observation that can be made is, for example

the consistent drop in positions of the item-based machine learning models in Movielens 20M indicating that the *correct* recommendations they provide are similar to those of the already available carousels, in particular the ItemKNN ones. This behavior depends on the dataset, with the models most negatively affected by the presence of an ItemKNN carousel being $P^3\alpha$ on ContentWise Impressions and $EASE^R$ on the Netflix Prize dataset.

These results indicate the importance of accounting how a set of recommendation lists complement each other and that this effect can change substantially the relative ranking of some algorithms compared to when they are evaluated individually, therefore leading to different conclusions on which is the best recommendation list to include.

7.2 Selecting a carousel layout

Table 5 compares the recommendation quality of the carousel layout selected with different approaches on interfaces with increasing number of carousels, from 2 to 8. The recommendation quality is measured with, and optimized for, NDCG at 10. The same experiment was also conducted optimizing N2DCG and produced consistent results. Although the maximum number of carousels reported here is 8, note that some content providers use much longer user interfaces, with Netflix showing more than 30 carousels.

As a general observation we can see how the solution space grows markedly and becomes unpractical to explore exhaustively even for rather small interfaces. For example, on Movielens 20M an interface of 4 carousels corresponds to $1.8 \cdot 10^3$ possible selections and $4.3 \cdot 10^4$ possible rankings. We limit the analysis to the exhaustive searches that could complete in a week of computation. Note that depending on the dataset this corresponds to very different numbers of layouts, from a minimum of $3.6 \cdot 10^2$ for Netflix Prize to a maximum of $9.4 \cdot 10^4$ for ContentWise Impressions. In this experiment, the default ranking adopted by the Exhaustive Selection method is to order the models according to their decreasing recommendation quality when evaluated individually.

Overall, the results indicate that the Exhaustive Selection achieves almost identical overall recommendation quality when compared to the Exhaustive Permutation in the limited number of cases where it is possible to use it. Indeed a difference may emerge for higher number of carousels but the growth of the search space makes such an analysis impractical. A more interesting analysis can be done with the proposed two greedy strategies. The Incremental Greedy method provides better results than the Individual Greedy, indicating again the importance of accounting for how the recommendation list complement each other. This difference is more marked on denser datasets, i.e., Netflix Prize, and for longer carousel layouts. The only dataset in which it is possible to apply the Exhaustive Selection up to 8 carousels is ContentWise Impressions, where we can see how the Incremental Greedy selects a layout of even better overall recommendation quality than the Exhaustive Selection. This should not be surprising since the exhaustive search strategies find the global optima, which makes them more prone to overfitting and hence exhibiting reduced generalizability. These results indicate that, although very simple, the Incremental Greedy strategy we described is indeed effective to account for how the different recommendation lists complement each other.

7.3 Comparing NDCG and N2DCG

The different results obtained by optimizing NDCG or N2DCG are not directly comparable in their absolute value, therefore are best visualized by comparing the corresponding optimal carousel layouts.

Table 6 compares the optimal layouts for 5 carousels obtained by all exhaustive searches and greedy strategies optimizing both NDCG and N2DCG. This is the longest interface for which it is possible to report results for all four strategies. By first comparing the exhaustive search strategies we can see that in both cases Exhaustive Search with Ranking selects the same models but there is always a couple that are swapped compared to the default Ranking strategy, indicating that a better ranking was found compared to the default one (i.e., decreasing order of individual accuracy). When optimizing the traditional single list NDCG both greedy strategies choose the same first two carousels the exhaustive strategies selected, but then the last three become different in both the ranking of the models and which ones are selected, i.e., Incremental Greedy selects $P^3\alpha$ while Individual Greedy does not and instead selects PureSVD, while $RP^3\beta$ and UserKNN CF are present in both but in different positions.

When optimizing the N2DCG the optimal layout is quite different, SLIM ElasticNet, that was consistently in first position when optimizing NDCG, is not even selected, while PureSVD is selected and put towards the end of the layout. The Incremental Greedy strategy selects the same first two carousels that were selected when optimizing NDCG, but the remaining three include different models with a different ranking. This highlights the intrinsic limitations of the simple Incremental Greedy strategy which, due to its definition, will not be able to account for the two dimensional structure of the interface when the very first carousel is selected and will therefore more likely find a suboptimal layout.

Similar observations can be made for the other two datasets and for a longer user interface. Table 7 shows the optimal layouts of 8 carousels selected for Netflix Prize and Movielens 20M. On Movielens 20M again the Incremental and Individual Greedy strategies are similar only in the first carousels but then result in quite different rankings and selected models, and the same holds for the ranking obtained by optimizing N2DCG. The Netflix dataset is instead an example of case where the Incremental Greedy is ineffective and selects the same layout by optimizing NDCG or N2DCG, indicating this is a scenario where developing a better strategy will be particularly important.

Overall this indicates that optimizing N2DCG, especially for longer page layouts, is a more difficult and nuanced problem, which opens new directions for future research, e.g., representing it as a quadratic optimization problem as done by (Ferrari Dacrema et al., 2021).

8 CONCLUSION

This paper proposes a new offline evaluation protocol for a carousel user interface, where the recommendation quality of a model is not measured independently but rather is put into the context of other recommendation lists being already available to the users. The experimental analysis shows that the relative ranking of the personalized algorithms changes when accounting for the presence of other carousels in the interface. This confirms previous observations that the correlations between recommendation lists have an important role to play and should be taken into account during offline evaluation as well. Results also show the impact of accounting for how the users navigate a user interface, with the traditional NDCG and the proposed N2DCG resulting in different optimal layouts when applying both exhaustive and greedy strategies. As future work, online studies should be conducted to measure how closely the offline carousel evaluation is able to represent the user behavior. For example, one could estimate appropriate factors to be used by the discount term of the N2DCG according to the layout structure and the user action type or explore different relevance functions that attribute partial relevance to duplicate items as well. Another direction is the development of recommendation models specifically tailored for the carousel scenario, e.g., that target specific items that are not accurately recommended by other algorithms (unpopular items), as

well as new efficient strategies to build an optimal carousel layout possibly personalized to communities of users or even single users. Ultimately, the carousel evaluation protocol opens new research directions by allowing researchers to conduct offline evaluations in this industrially relevant scenarios and open a wide number of research possibilities in studying how to combine the strength of various models and techniques to provide the user with ever more accurate and interesting recommendations.

CONFLICT OF INTEREST STATEMENT

The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

All authors contributed to the conception and design of the study. M.F.D. and N.F. wrote the first draft of the manuscript and performed the experiments. All authors contributed to manuscript revision, read, and approved the submitted version.

DATA AVAILABILITY STATEMENT

The datasets used in this study are publicly available and can be found in the referenced papers in Section 6.3.

REFERENCES

- Adomavicius, G. and Kwon, Y. (2012). Improving aggregate recommendation diversity using ranking-based techniques. *IEEE Trans. Knowl. Data Eng.* 24, 896–911. doi:10.1109/TKDE.2011.15
- Antenucci, S., Boglio, S., Chioso, E., Dervishaj, E., Kang, S., Scarlatti, T., et al. (2018). Artist-driven layering and user's behaviour impact on recommendations in a playlist continuation scenario. In *Proceedings of the ACM Recommender Systems Challenge, RecSys Challenge 2018, Vancouver, BC, Canada, October 2, 2018* (ACM), 4:1–4:6. doi:10.1145/3267471.3267475
- Bendada, W., Salha, G., and Bontempelli, T. (2020). Carousel personalization in music streaming apps with contextual bandits. In *RecSys 2020: Fourteenth ACM Conference on Recommender Systems, Virtual Event, Brazil, September 22–26, 2020*, eds. R. L. T. Santos, L. B. Marinho, E. M. Daly, L. Chen, K. Falk, N. Koenigstein, and E. S. de Moura (ACM), 420–425. doi:10.1145/3383313.3412217
- Bennett, J., Lanning, S., et al. (2007). The netflix prize. In *Proceedings of KDD cup and workshop* (New York), vol. 2007, 35
- Burges, C. J. C., Shaked, T., Renshaw, E., Lazier, A., Deeds, M., Hamilton, N., et al. (2005). Learning to rank using gradient descent. In *Machine Learning, Proceedings of the Twenty-Second International Conference (ICML 2005), Bonn, Germany, August 7–11, 2005*, eds. L. D. Raedt and S. Wrobel (ACM), vol. 119 of *ACM International Conference Proceeding Series*, 89–96. doi:10.1145/1102351.1102363
- Chierichetti, F., Kumar, R., and Raghavan, P. (2011). Optimizing two-dimensional search results presentation. In *Proceedings of the Forth International Conference on Web Search and Web Data Mining, WSDM 2011, Hong Kong, China, February 9–12, 2011*, eds. I. King, W. Nejdl, and H. Li (ACM), 257–266. doi:10.1145/1935826.1935873
- Cichocki, A. and Phan, A. H. (2009). Fast local algorithms for large scale nonnegative matrix and tensor factorizations. *IEICE Trans. Fundam. Electron. Commun. Comput. Sci.* 92-A, 708–721. doi:10.1587/transfun.E92.A.708

- Cooper, C., Lee, S., Radzik, T., and Siantos, Y. (2014). Random walks in recommender systems: exact computation and simulations. In *23rd International World Wide Web Conference, WWW '14, Seoul, Republic of Korea, April 7-11, 2014, Companion Volume*, eds. C. Chung, A. Z. Broder, K. Shim, and T. Suel (ACM), 811–816. doi:10.1145/2567948.2579244
- Craswell, N., Zoeter, O., Taylor, M. J., and Ramsey, B. (2008). An experimental comparison of click position-bias models. In *Proceedings of the International Conference on Web Search and Web Data Mining, WSDM 2008, Palo Alto, California, USA, February 11-12, 2008*, eds. M. Najork, A. Z. Broder, and S. Chakrabarti (ACM), 87–94. doi:10.1145/1341531.1341545
- Cremonesi, P., Koren, Y., and Turrin, R. (2010). Performance of recommender algorithms on top-n recommendation tasks. In *Proceedings of the 2010 ACM Conference on Recommender Systems, RecSys 2010, Barcelona, Spain, September 26-30, 2010*, eds. X. Amatriain, M. Torrens, P. Resnick, and M. Zanker (ACM), 39–46. doi:10.1145/1864708.1864721
- Ding, W., Govindaraj, D., and Vishwanathan, S. V. N. (2019). Whole page optimization with global constraints. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining, KDD 2019, Anchorage, AK, USA, August 4-8, 2019*, eds. A. Teredesai, V. Kumar, Y. Li, R. Rosales, E. Terzi, and G. Karypis (ACM), 3153–3161. doi:10.1145/3292500.3330675
- Elahi, E. and Chandrashekar, A. (2020). Learning representations of hierarchical slates in collaborative filtering. In *RecSys 2020: Fourteenth ACM Conference on Recommender Systems, Virtual Event, Brazil, September 22-26, 2020*, eds. R. L. T. Santos, L. B. Marinho, E. M. Daly, L. Chen, K. Falk, N. Koenigstein, and E. S. de Moura (ACM), 703–707. doi:10.1145/3383313.3418484
- Felicioni, N., Ferrari Dacrema, M., and Cremonesi, P. (2021a). A methodology for the offline evaluation of recommender systems in a user interface with multiple carousels. In *Adjunct Publication of the 29th ACM Conference on User Modeling, Adaptation and Personalization, UMAP 2021, Utrecht, The Netherlands, June 21-25, 2021*, eds. J. Masthoff, E. Herder, N. Tintarev, and M. Tkalcic (ACM), 10–15. doi:10.1145/3450614.3461680
- Felicioni, N., Ferrari Dacrema, M., Pérez Maurera, F. B., and Cremonesi, P. (2021b). Measuring the ranking quality of recommendations in a two-dimensional carousel setting. In *Proceedings of the 11th Italian Information Retrieval Workshop 2021, Bari, Italy, September 13-15, 2021*, eds. V. W. Anelli, T. D. Noia, N. Ferro, and F. Narducci (CEUR-WS.org), vol. 2947 of *CEUR Workshop Proceedings*
- Ferrari Dacrema, M. (2021). Demonstrating the equivalence of list based and aggregate metrics to measure the diversity of recommendations (student abstract). In *Thirty-Fifth AAAI Conference on Artificial Intelligence, AAAI 2021, Thirty-Third Conference on Innovative Applications of Artificial Intelligence, IAAI 2021, The Eleventh Symposium on Educational Advances in Artificial Intelligence, EAAI 2021, Virtual Event, February 2-9, 2021* (AAAI Press), 15779–15780
- Ferrari Dacrema, M., Boglio, S., Cremonesi, P., and Jannach, D. (2021). A troubling analysis of reproducibility and progress in recommender systems research. *ACM Trans. Inf. Syst.* 39. doi:10.1145/3434185
- Ferrari Dacrema, M., Cremonesi, P., and Jannach, D. (2019). Are we really making much progress? A worrying analysis of recent neural recommendation approaches. In *Proceedings of the 13th ACM Conference on Recommender Systems, RecSys 2019, Copenhagen, Denmark, September 16-20, 2019*, eds. T. Bogers, A. Said, P. Brusilovsky, and D. Tikk (ACM), 101–109. doi:10.1145/3298689.3347058
- Ferrari Dacrema, M., Felicioni, N., and Cremonesi, P. (2021). Optimizing the selection of recommendation carousels with quantum computing. In *RecSys '21: Fifteenth ACM Conference on Recommender Systems, Amsterdam, The Netherlands, 27 September 2021 - 1 October 2021*, eds. H. J. C. Pampín, M. A. Larson,

- 758 M. C. Willemsen, J. A. Konstan, J. J. McAuley, J. Garcia-Gathright, B. Huurnink, and E. Oldridge
759 (ACM), 691–696. doi:10.1145/3460231.3478853
- 760 Gruson, A., Chandar, P., Charbuillet, C., McInerney, J., Hansen, S., Tardieu, D., et al. (2019). Offline
761 evaluation to make decisions about playlist recommendation algorithms. In *Proceedings of the Twelfth*
762 *ACM International Conference on Web Search and Data Mining, WSDM 2019, Melbourne, VIC,*
763 *Australia, February 11–15, 2019*, eds. J. S. Culpepper, A. Moffat, P. N. Bennett, and K. Lerman (ACM),
764 420–428. doi:10.1145/3289600.3291027
- 765 Harper, F. M. and Konstan, J. A. (2016). The movielens datasets: History and context. *ACM Trans. Interact.*
766 *Intell. Syst.* 5, 19:1–19:19. doi:10.1145/2827872
- 767 Herlocker, J. L., Konstan, J. A., Terveen, L. G., and Riedl, J. (2004). Evaluating collaborative filtering
768 recommender systems. *ACM Trans. Inf. Syst.* 22, 5–53. doi:10.1145/963770.963772
- 769 Hu, Y., Koren, Y., and Volinsky, C. (2008). Collaborative filtering for implicit feedback datasets. In
770 *Proceedings of the 8th IEEE International Conference on Data Mining (ICDM 2008), December 15–19,*
771 *2008, Pisa, Italy* (IEEE Computer Society), 263–272. doi:10.1109/ICDM.2008.22
- 772 Järvelin, K. and Kekäläinen, J. (2000). IR evaluation methods for retrieving highly relevant documents. In
773 *SIGIR 2000: Proceedings of the 23rd Annual International ACM SIGIR Conference on Research and*
774 *Development in Information Retrieval, July 24–28, 2000, Athens, Greece*, eds. E. J. Yannakoudakis, N. J.
775 Belkin, P. Ingwersen, and M. Leong (ACM), 41–48. doi:10.1145/345508.345545
- 776 Järvelin, K. and Kekäläinen, J. (2002). Cumulated gain-based evaluation of IR techniques. *ACM Trans. Inf.*
777 *Syst.* 20, 422–446. doi:10.1145/582415.582418
- 778 Kammerer, Y. and Gerjets, P. (2010). How the interface design influences users' spontaneous
779 trustworthiness evaluations of web search results: comparing a list and a grid interface. In *Proceedings*
780 *of the 2010 Symposium on Eye-Tracking Research & Applications, ETRA 2010, Austin, Texas, USA,*
781 *March 22–24, 2010*, eds. C. H. Morimoto, H. O. Istance, A. Hyrskykari, and Q. Ji (ACM), 299–306.
782 doi:10.1145/1743666.1743736
- 783 Kanoulas, E. and Aslam, J. A. (2009). Empirical justification of the gain and discount function for ndcg.
784 In *Proceedings of the 18th ACM Conference on Information and Knowledge Management, CIKM 2009,*
785 *Hong Kong, China, November 2–6, 2009*, eds. D. W. Cheung, I. Song, W. W. Chu, X. Hu, and J. J. Lin
786 (ACM), 611–620. doi:10.1145/1645953.1646032
- 787 Mobasher, B., Jin, X., and Zhou, Y. (2003). Semantically enhanced collaborative filtering on the web. In
788 *Web Mining: From Web to Semantic Web, First European Web Mining Forum, EMWF 2003, Cavtat-*
789 *Dubrovnik, Croatia, September 22, 2003, Revised Selected and Invited Papers*, eds. B. Berendt, A. Hotho,
790 D. Mladenic, M. van Someren, M. Spiliopoulou, and G. Stumme (Springer), vol. 3209 of *Lecture Notes*
791 *in Computer Science*, 57–76. doi:10.1007/978-3-540-30123-3_4
- 792 Ning, X. and Karypis, G. (2011). SLIM: Sparse linear methods for top-n recommender systems. In
793 *Proceedings of the 11th IEEE International Conference on Data Mining (ICDM '11)*. 497–506
- 794 Paudel, B., Christoffel, F., Newell, C., and Bernstein, A. (2017). Updatable, accurate, diverse, and
795 scalable recommendations for interactive applications. *ACM Trans. Interact. Intell. Syst.* 7, 1:1–1:34.
796 doi:10.1145/2955101
- 797 Pérez Maurera, F. B., Ferrari Dacrema, M., Saule, L., Scriminaci, M., and Cremonesi, P. (2020).
798 Contentwise impressions: An industrial dataset with impressions included. In *CIKM '20: The 29th ACM*
799 *International Conference on Information and Knowledge Management, Virtual Event, Ireland, October*
800 *19–23, 2020*, eds. M. d'Aquin, S. Dietze, C. Hauff, E. Curry, and P. Cudré-Mauroux (ACM), 3093–3100.
801 doi:10.1145/3340531.3412774

- 802 Rendle, S., Freudenthaler, C., Gantner, Z., and Schmidt-Thieme, L. (2009). BPR: bayesian personalized
803 ranking from implicit feedback. In *UAI 2009, Proceedings of the Twenty-Fifth Conference on Uncertainty*
804 *in Artificial Intelligence, Montreal, QC, Canada, June 18-21, 2009*, eds. J. A. Bilmes and A. Y. Ng
805 (AUAI Press), 452–461
- 806 Sanderson, M. and Croft, W. B. (2012). The history of information retrieval research. *Proc. IEEE* 100,
807 1444–1451. doi:10.1109/JPROC.2012.2189916
- 808 Sarwar, B. M., Karypis, G., Konstan, J. A., and Riedl, J. (2001). Item-based collaborative filtering
809 recommendation algorithms. In *Proceedings of the Tenth International World Wide Web Conference,*
810 *WWW 10, Hong Kong, China, May 1-5, 2001*, eds. V. Y. Shen, N. Saito, M. R. Lyu, and M. E. Zurko
811 (ACM), 285–295. doi:10.1145/371920.372071
- 812 Steck, H. (2019). Embarrassingly shallow autoencoders for sparse data. In *The World Wide Web Conference,*
813 *WWW 2019, San Francisco, CA, USA, May 13-17, 2019*, eds. L. Liu, R. W. White, A. Mantrach,
814 F. Silvestri, J. J. McAuley, R. Baeza-Yates, and L. Zia (ACM), 3251–3257. doi:10.1145/3308558.
815 3313710
- 816 Wu, C., Alvino, C. V., Smola, A. J., and Basilico, J. (2016). Using navigation to improve recommendations
817 in real-time. In *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA,*
818 *USA, September 15-19, 2016*, eds. S. Sen, W. Geyer, J. Freyne, and P. Castells (ACM), 341–348.
819 doi:10.1145/2959100.2959174
- 820 Zhao, Q., Chang, S., Harper, F. M., and Konstan, J. A. (2016). Gaze prediction for recommender systems. In
821 *Proceedings of the 10th ACM Conference on Recommender Systems, Boston, MA, USA, September 15-19,*
822 *2016*, eds. S. Sen, W. Geyer, J. Freyne, and P. Castells (ACM), 131–138. doi:10.1145/2959100.2959150
- 823 Zhou, K., Zha, H., Chang, Y., and Xue, G. (2014). Learning the gain values and discount factors of
824 discounted cumulative gains. *IEEE Trans. Knowl. Data Eng.* 26, 391–404. doi:10.1109/TKDE.2012.252
- 825 Zhou, T., Kuscsik, Z., Liu, J.-G., Medo, M., Wakeling, J. R., and Zhang, Y.-C. (2010). Solving the apparent
826 diversity-accuracy dilemma of recommender systems. *National Academy of Sciences* 107, 4511–4515

FIGURE CAPTIONS

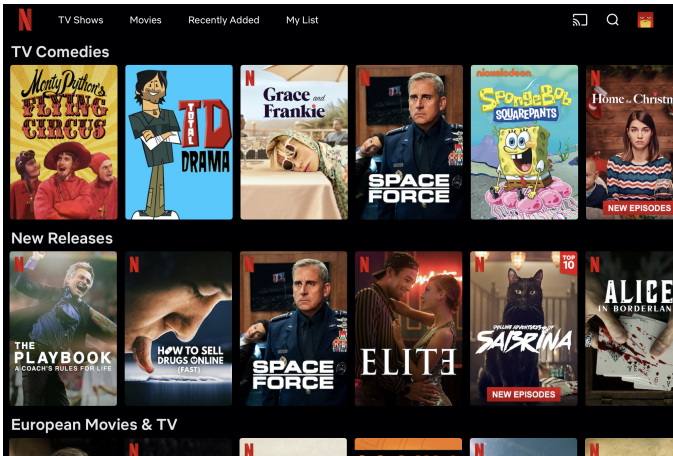


Figure 1. The Netflix homepage, an example of carousel user interface in the multimedia streaming domain.

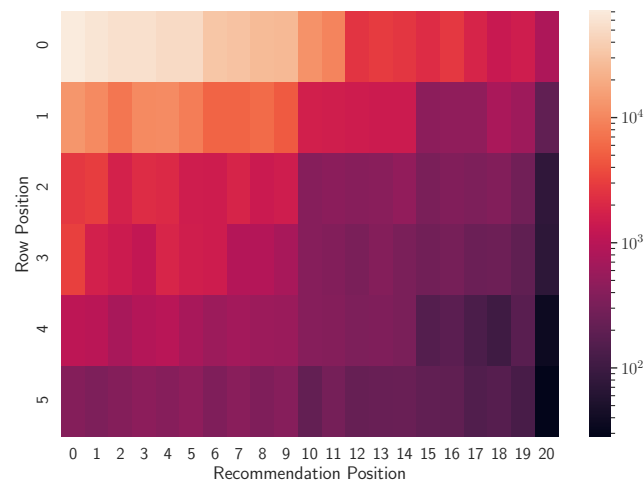


Figure 2. Heatmap of the number of interactions per position on the screen, taken from the dataset presented in (Pérez Maurera et al., 2020).

TABLES

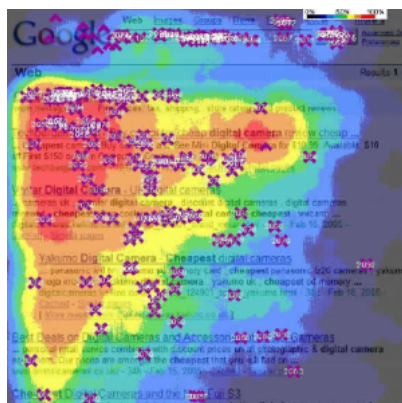


Figure 3. Heatmap from real usage data of how users navigate a simple search result page, the attention is focused on a top-left *golden triangle* and progressively decreases moving to the right and the bottom (Chierichetti et al., 2011).

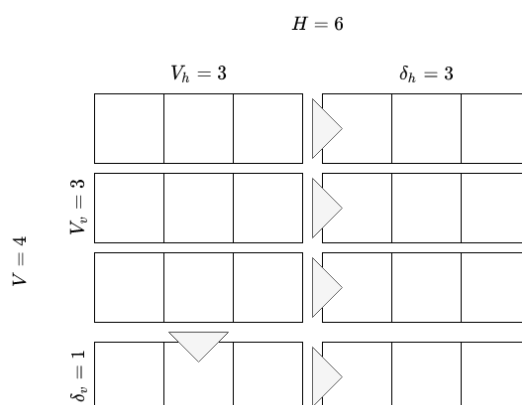


Figure 4. An example interface with 4 carousels of 6 elements each, where 3 carousels, with 3 items each, are visible. A horizontal swipe reveals 3 items, while a vertical swipe reveals one additional carousel.

Table 1. Summary of the notation used for the carousel interface.

Symbol	Description
H	Horizontal dimension of the interface, i.e., number of columns. Corresponds to the length of each recommendation list.
V	Vertical dimension of the interface, i.e., number of rows. Corresponds to the number of carousels.
M	Total number of recommendation lists that can be used for the carousel layout, $M \geq V$.
l	Total number of recommendations provided to the user. When evaluating a single list $l = H$. When evaluating a carousel layout $l = HV$.
V_h	Number of columns of the interface that are immediately visible to the user, $V_h \leq H$
V_v	Number of rows of the interface that are immediately visible to the user, $V_v \leq V$
δ_h	Number of columns that are revealed after a user action, $\delta_h \leq V_h$
δ_v	Number of rows that are revealed after a user action, $\delta_v \leq V_v$

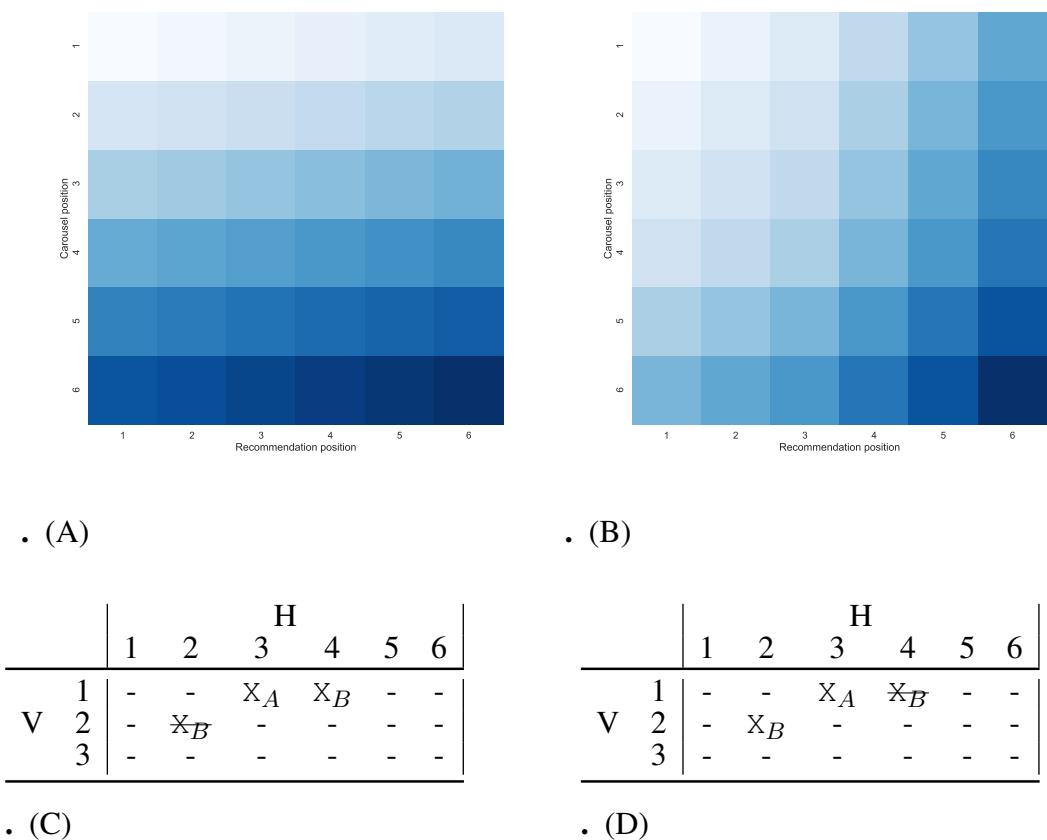


Figure 5. A visual comparison of the two-dimensional discount function under different assumptions: Figure (A) represents carousels concatenated in a single list. Figure (B) represents two-dimensional discount which accounts for *golden triangle* and number of user actions. Figure (C) and (D) show the impact of the discount function on the relevance of a correct recommendation. Consider a small user interface with 3 carousels, each with 6 elements. Correct recommendations are represented with “ X_n ” with n the item identifier, while incorrect ones by “-”. Item B is recommended correctly both by the first and second recommendation list. Only one of the two will be considered relevant, depending on the discount, the other is shown as crossed out. The single list discount shown in Figure (A) will lead to consider relevant the recommendation of item B in the *first* carousel, while the two-dimensional discount shown in Figure (B) will lead to consider relevant the recommendation in the *second* carousel.

		H					
		1	2	3	4	5	6
V	1	-	X	-	-	-	-
	2	-	-	X	-	-	-
	3	-	X	-	-	-	-

• (A)

		H					
		1	2	3	4	5	6
V	1	-	X	-	-	-	-
	2	-	-	X	-	-	-
	3	X	X	-	-	-	-

• (B)

		H					
		1	2	3	4	5	6
V	1	-	-	X	X	-	-
	2	-	X	-	-	-	-
	3	-	-	-	-	-	-

• (C)

		H					
		1	2	3	4	5	6
V	1	-	X	-	-	-	-
	2	-	-	X	X	-	-
	3	-	-	-	-	-	-

• (D)

Figure 6. A small user interface with 3 carousels, each with 6 elements on which the behavior of DCG and 2DCG is compared. Correct recommendations are represented with “X” while incorrect ones by “-”. The vertical line between column 3 and 4 represents the limit of the area initially visible to the user, i.e., $V_h = 3$. Figure (A) and (B) represent the same layout in which a correct recommendation in position (3, 1) is added. Figure (C) and (D) represent two different layouts created based on the same three recommendation lists, with DCG preferring layout (C) and 2DCG preferring layout (D) due to the effect of the user action penalty on the correct recommendation in column 4.

Table 2. Comparison of the NDCG at 10 and overall model ranking for Movielens 20M. Each model is evaluated both individually (single-carousel) and as the last recommendation list in a multi-carousel interface of increasing complexity. The NDCG is computed with the single-list discount (concatenating all carousel lists). Higher ranks indicate better recommendation quality. The rank of models that are already used as carousels is removed. Δ Rank is the difference between the rank when evaluated individually and the rank when evaluated in the corresponding carousel layout, e.g., a negative Δ Rank indicates the model is in a worse ranking position.

	Individual		TopPop		Carousel Layout TopPop ItemKNN CF		TopPop ItemKNN CF ItemKNN CBF	
	NDCG	Rank	NDCG	Δ Rank	NDCG	Δ Rank	NDCG	Δ Rank
TopPopular	0.1058	15	0.0953	-	0.1607	-	0.1713	-
ItemKNN CF	0.2216	11	0.1653	0	0.1607	-	0.1713	-
ItemKNN CBF	0.1202	14	0.1271	0	0.1734	3	0.1713	-
GlobalEffects	0.0478	16	0.1041	0	0.1659	0	0.1747	0
UserKNN CF	0.3088	2	0.1917	0	0.1958	1	0.1978	1
$P^3\alpha$	0.1968	12	0.1393	0	0.1695	0	0.1778	1
$RP^3\beta$	0.2422	8	0.1667	0	0.1743	-1	0.1812	-1
IALS	0.2702	5	0.1868	2	0.1932	2	0.1961	2
MF BPR	0.1645	13	0.1352	0	0.1689	0	0.1776	1
MF FunkSVD	0.2519	7	0.1790	2	0.1922	3	0.1952	3
PureSVD	0.2657	6	0.1765	0	0.1832	0	0.1872	0
NMF	0.2288	9	0.1662	0	0.1819	1	0.1871	2
EASE ^R	0.2740	4	0.1756	-3	0.1823	-3	0.1861	-4
SLIM ElasticNet	0.3109	1	0.1920	0	0.1951	-1	0.1964	-1
SLIM BPR	0.2792	3	0.1811	-1	0.1882	-2	0.1918	-2
ItemKNN CFCBF	0.2264	10	0.1662	0	0.1667	-3	0.1759	-2

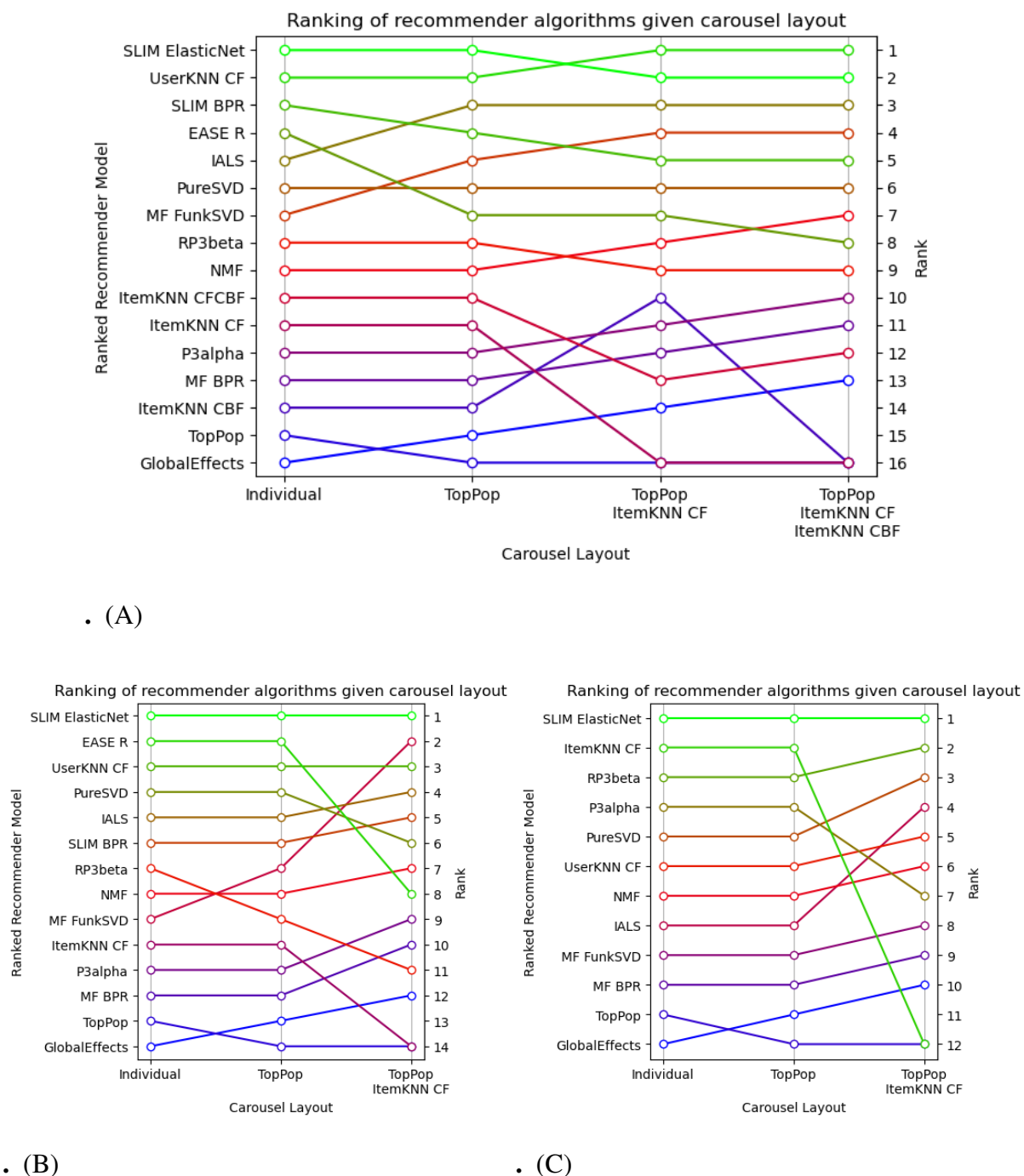


Figure 7. Visualization of how the ranking of several recommendation models changes when they are evaluated independently or as the last recommendation list in a carousel interface of increasing complexity. Figure (A) Movielens 20M, Figure (B) Netflix Prize and Figure (C) ContentWise Impressions Highest ranked models are the best performing according to NDCG at 10.

Table 3. Comparison of the NDCG at 10 value and overall model ranking for ContentWise Impressions. Each model is evaluated both individually (single-carousel) and as the last recommendation list in a multi-carousel interface of increasing complexity. The NDCG is computed with the single-list discount (concatenating all carousel lists). Higher ranks indicate better recommendation quality. The rank of models that are already used as carousels is removed. Δ Rank is the difference between the rank when evaluated individually and the rank when evaluated in the corresponding carousel layout, e.g., a negative Δ Rank indicates the model is in a worse ranking position.

	Individual		Carousel Layout			
			TopPop		TopPop ItemKNN CF	
	NDCG	Rank	NDCG	Δ Rank	NDCG	Δ Rank
TopPopular	0.0708	11	0.0617	-	0.2448	-
ItemKNN CF	0.5328	2	0.2557	0	0.2448	-
GlobalEffects	0.0000	12	0.0617	0	0.2448	0
UserKNN CF	0.3667	6	0.1915	0	0.2688	0
$P^3\alpha$	0.4183	4	0.2095	0	0.2686	-4
$RP^3\beta$	0.5016	3	0.2479	0	0.2745	0
IALS	0.3369	8	0.1754	0	0.2696	3
MF BPR	0.1584	10	0.1118	0	0.2554	0
MF FunkSVD	0.2324	9	0.1395	0	0.2624	0
PureSVD	0.3699	5	0.1935	0	0.2712	1
NMF	0.3428	7	0.1798	0	0.2686	0
SLIM ElasticNet	0.5548	1	0.2651	0	0.2769	0

Table 4. Comparison of the NDCG at 10 value and overall model ranking for NetflixPrize. Each model is evaluated both individually (single-carousel) and as the last recommendation list in a multi-carousel interface of increasing complexity. The NDCG is computed with the single-list discount (concatenating all carousel lists). Higher ranks indicate better recommendation quality. The rank of models that are already used as carousels is removed. Δ Rank is the difference between the rank when evaluated individually and the rank when evaluated in the corresponding carousel layout, e.g., a negative Δ Rank indicates the model is in a worse ranking position.

	Individual		Carousel Layout			
			TopPop		TopPop ItemKNN CF	
	NDCG	Rank	NDCG	Δ Rank	NDCG	Δ Rank
TopPopular	0.0799	13	0.0678	-	0.1261	-
ItemKNN CF	0.2060	10	0.1335	0	0.1261	-
GlobalEffects	0.0159	14	0.0733	0	0.1282	0
UserKNN CF	0.2581	3	0.1516	0	0.1573	0
$P^3\alpha$	0.1810	11	0.1194	0	0.1478	1
$RP^3\beta$	0.2209	7	0.1345	-2	0.1422	-4
IALS	0.2380	5	0.1497	0	0.1566	1
MF BPR	0.1656	12	0.1172	0	0.1426	1
MF FunkSVD	0.2077	9	0.1446	2	0.1639	7
PureSVD	0.2508	4	0.1515	0	0.1552	-2
NMF	0.2192	8	0.1434	0	0.1546	1
$EASE^R$	0.2619	2	0.1520	0	0.1534	-6
SLIM ElasticNet	0.2913	1	0.1662	0	0.1669	0
SLIM BPR	0.2353	6	0.1467	0	0.1561	1

Table 5. Comparison of the NDCG at 10 of the layout selected according to the two proposed greedy strategies and compared to the exhaustive search of all possible model selections and rankings. The page layout contains from 2 to 8 carousels. Results for exhaustive searches requiring more than a week of computation are missing.

		Number of Carousels						
		2	3	4	5	6	7	8
Movielens 20M ($M = 16$)								
Solution Space Size	Number of Rankings	$2.4 \cdot 10^2$	$3.3 \cdot 10^3$	$4.3 \cdot 10^4$	$5.2 \cdot 10^5$	$5.7 \cdot 10^6$	$5.7 \cdot 10^7$	$5.1 \cdot 10^8$
	Number of Selections	$1.2 \cdot 10^2$	$5.6 \cdot 10^2$	$1.8 \cdot 10^3$	$4.3 \cdot 10^3$	$8.0 \cdot 10^3$	$1.1 \cdot 10^4$	$1.2 \cdot 10^4$
Selection Method	Exhaustive Selection and Ranking	0.3168	0.3220	-	-	-	-	-
	Exhaustive Selection default Ranking	0.3168	0.3220	0.3266	0.3296	-	-	-
	Incremental Greedy	0.3168	0.3220	0.3266	0.3297	0.3320	0.3345	0.3363
	Individual Greedy	0.3168	0.3180	0.3190	0.3241	0.3250	0.3295	0.3309
Netflix Prize ($M = 14$)								
Solution Space Size	Number of Rankings	$1.8 \cdot 10^2$	$2.1 \cdot 10^3$	$2.4 \cdot 10^4$	$2.4 \cdot 10^5$	$2.1 \cdot 10^6$	$1.7 \cdot 10^7$	$1.2 \cdot 10^8$
	Number of Selections	$9.1 \cdot 10^1$	$3.6 \cdot 10^2$	$1.0 \cdot 10^3$	$2.0 \cdot 10^3$	$3.0 \cdot 10^3$	$3.4 \cdot 10^3$	$3.0 \cdot 10^3$
Selection Method	Exhaustive Selection and Ranking	0.2781	-	-	-	-	-	-
	Exhaustive Selection default Ranking	0.2781	0.2809	-	-	-	-	-
	Incremental Greedy	0.2781	0.2812	0.2830	0.2865	0.2891	0.2912	0.2931
	Individual Greedy	0.2692	0.2701	0.2689	0.2714	0.2740	0.2761	0.2777
ContentWise Impressions ($M = 12$)								
Solution Space Size	Number of Rankings	$1.3 \cdot 10^2$	$1.3 \cdot 10^3$	$1.1 \cdot 10^4$	$9.5 \cdot 10^4$	$6.6 \cdot 10^5$	$3.9 \cdot 10^6$	$2.0 \cdot 10^7$
	Number of Selections	$6.6 \cdot 10^1$	$2.2 \cdot 10^2$	$4.9 \cdot 10^2$	$7.9 \cdot 10^2$	$9.2 \cdot 10^2$	$7.9 \cdot 10^2$	$4.9 \cdot 10^2$
Selection Method	Exhaustive Selection and Ranking	0.5162	0.5099	0.5133	0.5153	-	-	-
	Exhaustive Selection default Ranking	0.5162	0.5099	0.5132	0.5152	0.5178	0.5199	0.5213
	Incremental Greedy	0.5162	0.5099	0.5133	0.5153	0.5179	0.5201	0.5215
	Individual Greedy	0.5162	0.5098	0.5123	0.5149	0.5172	0.5183	0.5190

Table 6. Layouts of 5 carousels selected for the ContentWise Impressions dataset selected according to different strategies and optimizing both the single list NDCG and the proposed two-dimensional N2DCG. The Individual Greedy layout for N2DCG is missing since that method is based on the traditional single list evaluation.

Optimized Metric	Exhaustive Selection and Ranking	Exhaustive Selection default Ranking	Incremental Greedy	Individual Greedy
NDCG	SLIM ElasticNet	SLIM ElasticNet	SLIM ElasticNet	SLIM ElasticNet
	ItemKNN CF	ItemKNN CF	ItemKNN CF	ItemKNN CF
	UserKNN CF	RP ³ β	RP ³ β	UserKNN CF
	RP ³ β	UserKNN CF	P ³ α	RP ³ β
	FunkSVD	FunkSVD	UserKNN CF	FunkSVD
N2DCG	ItemKNN CF	ItemKNN CF	SLIMElasticNet	-
	P ³ α	RP ³ β	ItemKNN CF	-
	RP ³ β	P ³ α	UserKNN CF	-
	PureSVD	PureSVD	RP ³ β	-
	FunkSVD	FunkSVD	FunkSVD	-

Table 7. Layouts of 8 carousels selected for the Movielens 20M and Netflix Prize datasets selected according to different strategies and optimizing both the single list NDCG and the proposed two-dimensional N2DCG. The Individual Greedy layout for N2DCG is missing since that method is based on the traditional single list evaluation.

Movielens 20M			Netflix Prize		
NDCG		N2DCG	NDCG		N2DCG
Incremental Greedy	Individual Greedy	Incremental Greedy	Incremental Greedy	Individual Greedy	Incremental Greedy
UserKNN CF	UserKNN CF	UserKNN CF	SLIM ElasticNet	SLIM ElasticNet	SLIM ElasticNet
SLIM ElasticNet	SLIM ElasticNet	SLIM ElasticNet	FunkSVD	EASE ^R	FunkSVD
FunkSVD	SLIM BPR	FunkSVD	UserKNN CF	UserKNN CF	UserKNN CF
IALS	EASE ^R	IALS	MF BPR	PureSVD	MF BPR
MF BPR	IALS	MF BPR	IALS	SLIM BPR	IALS
NMF	PureSVD	ItemKNN CBF	P ³ _α	IALS	P ³ _α
ItemKNN CBF	FunkSVD	NMF	NMF	RP ³ _β	NMF
ItemKNN CF	RP ³ _β	ItemKNN CF	ItemKNN CF	ItemKNN CF	ItemKNN CF

1 OPTIMAL HYPERPARAMETERS

This section reports the optimal hyperparameters found for all recommendation models and datasets. The hyperparameter ranges and distribution used during the optimization are the same used and described in (Ferrari Dacrema et al., 2021). Table 8 refers to the collaborative KNN models, 9 to the content based and hybrid KNN models and 10 to the machine learning and graph based models.

Table 8. Hyperparameter values for collaborative KNN recommender algorithms on all datasets.

Algorithm	Hyperparameter	ContentWise Impressions	Netflix Prize	Movielens 20M
UserKNN CF	topK	297	1000	835
	shrink	102	0	0
	similarity	cosine	cosine	cosine
	normalize	True	True	True
	feature weighting	TF-IDF	none	none
ItemKNN CF	topK	5	72	214
	shrink	0	0	982
	similarity	cosine	cosine	cosine
	normalize	True	True	True
	feature weighting	none	TF-IDF	none

Table 9. Hyperparameter values for content based and hybrid KNN recommender algorithms. Only Movielens 20M contains item features, the other datasets are therefore omitted.

Algorithm	Hyperparameter	Movielens 20M
ItemKNN CBF	topK	1000
	shrink	1000
	similarity	cosine
	normalize	True
	feature weighting	TF-IDF
ItemKNN CF CBF	topK	253
	shrink	203
	similarity	cosine
	normalize	True
	feature weighting	BM25
	ICM weight	0.1520

Table 10. Hyperparameter values for machine learning and graph based recommender algorithms on all datasets. Results for EASE^R and SLIM BPR are missing for the ContentWise Impressions dataset because the algorithms required more than the available 64GB of RAM.

Algorithm	Hyperparameter	ContentWise Impressions	Netflix Prize	Movielens 20M
P ³ _α	topK	146	5	1000
	alpha	1.2748	0.9150	1.0915
	normalize similarity	True	False	True
RP ³ _β	topK	5	173	245
	alpha	0.0000	0.7040	0.5775
	beta	0.3487	0.6578	0.5625
	normalize similarity	True	True	True
SLIM ElasticNet	topK	1000	503	588
	l1 ratio	1.14E-02	1.33E-02	2.17E-01
	alpha	0.0010	0.0010	0.0023
MF BPR	sgd mode	adagrad	adagrad	adagrad
	epochs	1415	1500	1065
	num factors	200	200	158
	batch size	2	1	256
	positive reg	2.32E-04	1.00E-05	1.00E-05
	negative reg	1.50E-04	1.00E-05	7.03E-05
	learning rate	5.86E-02	2.96E-02	2.16E-02
MF FunkSVD	sgd mode	adam	adam	adagrad
	epochs	500	280	160
	use bias	False	False	False
	batch size	32	1024	2
	num factors	184	186	127
	item reg	3.92E-04	5.98E-04	1.47E-05
	user reg	2.84E-03	9.44E-04	3.18E-03
	learning rate	1.62E-04	1.10E-03	7.11E-02
	negative quota	0.0133	0.1198	0.0266
PureSVD	num factors	350	47	29
NMF	num factors	324	79	79
	solver	mult. update	mult. update	mult. update
	init type	nndsvda	random	random
	beta loss	frobenius	frobenius	kullback-leibler
IALS	num factors	200	63	60
	epochs	5	65	60
	confidence scaling	linear	linear	linear
	alpha	0.7942	0.2199	0.7748
	epsilon	0.1012	0.0010	0.1358
	reg	1.00E-02	1.00E-02	1.00E-02
EASE ^R	l2 norm	-	2.12E+06	7.52E+05
SLIM BPR	topK	-	36	685
	epochs	-	595	970
	symmetric	-	True	True
	sgd mode	-	sgd	sgd
	lambda i	-	2.07E-03	1.00E-05
	lambda j	-	9.20E-03	1.00E-05
	learning rate	-	9.22E-04	1.56E-03