

Quantifying Exploration Preference for E-Commerce Recommendation

Amy B.Z. Zhang^{1,*}, Siyun Wang², Raphael Louca², Karl Ni² and Diane Hu²

¹Cornell University, 2 W Loop Rd, New York, NY 10044, USA

²Etsy, Inc., 117 Adams St, Brooklyn, NY 11201, USA

Abstract

Recommendation systems are widely used in e-commerce setting to help users find the most relevant items. However, what is relevant for a user changes dynamically. For next-item recommendation, similarity to recently interacted items is desirable in some cases but not others: consider for example the behavior during comparison shopping in contrast to after a purchase. Such shifting user preference regarding exploration is not well captured by existing concepts, much less taken into account in recommendations. In this paper, we offer definitions to quantify user exploration preference and how it trends over time, based on spread of recently interacted items in embedding space. The soundness of the concepts are illustrated with mathematical properties as well as analysis of platform data. We further demonstrate flexibility and potential by attaching simple modules to well-known baseline algorithms in two separate use cases, comparing performances for three e-commerce datasets. Source code can be found at <https://github.com/bz275/DispPred>.

Keywords

sequential recommendation, dynamic preference, user profiling, embedding model, adaptive module, evaluation metric

1. Introduction

It is well recognized that user satisfaction with next-item recommendation depends on far more than similarity to past interactions [1, 2, 3]. User preferences shift dynamically with context and intent, making accounting for such changes important for relevant recommendations.

One aspect of this preference is with regard to the preference for exploration. Consider the example stream of user activity illustrated by Figure 1, where the user exhibited a change in attitude towards unfamiliar items. It is intuitive to suppose a user engaging in comparison shopping is interested in seeing similar alternatives to previous visited items, but prefers novelty after the purchase is complete. However there lacks quantitative measure that can reflect such trends, so the focus of this work is to quantify a user's exploration preference in a dynamic fashion.

Specifically, we consider the context of sequential recommendation in e-commerce, and quantify the *affinity to variation* a user demonstrates at a given time based on the spread of


ORSUM@ACM RecSys 2023: 6th Workshop on Online Recommender Systems and User Modeling, jointly with the 17th ACM Conference on Recommender Systems, September 19th, 2023, Singapore

*Work conducted while employed at Etsy, Inc.

✉ bz275@cornell.edu (A. B.Z. Zhang); swang@etsy.com (S. Wang); rlouca@etsy.com (R. Louca); kni@etsy.com (K. Ni); dhu@etsy.com (D. Hu)



© 2023 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

 CEUR Workshop Proceedings (CEUR-WS.org)

their recent activities. With access to item embeddings that reflect item similarity/ relation, the "variance" of the past L items quantifies the *Dispersion* at that timestep for look-back window L . A smaller value indicates that the vectors are more concentrated, suggesting focused activities, while a larger value indicates greater dispersion, implying more scattered browsing. As the look-back window slides forward with each additional interaction, the relative change, *DispAdd*, reflects the direction and magnitude of how the affinity is trending. A positive value suggests a trend towards increased variation, and the magnitude against the bounded scale indicates the strength of this trend.



Figure 1: Item images corresponding to a sequence of interactions for a e-commerce customer, which starts with narrow focus leading to a purchase, expands to more categories, then narrows again for a second purchase.

To show how the proposed concepts provide valuable insights into user preference changes in a real-world setting, we use both publicly available e-commerce datasets and data collected from live traffic on Etsy, a two-sided online marketplace notable for unique items. We offer data analysis results to illustrate conceptually, and also provide proof-of-concept experiments for two separate tasks, next-item recommendation and purchase prediction. Specifically, we show that incorporating *DispAdd* improves accuracy (Hit Rate and NDCG) in the former (with lower *Dispersion* deviation), boosts both ROC-AUC and Precision-Recall AUC in the latter, even with the module added using the most basic structure. The embeddings we use are generated using the Universal Sentence Encoder [4], which is more general, as well as SASRec [5], which is task specific. All experiments are conducted with SASRec as the base engine; while later methods (e.g. [6], [7]) have competitive performances, SASRec has shown robust SOTA performance and is representative of the basic transformer infrastructure most commonly used in practical systems. The adaptive module is added using the most basic structure, demonstrating the potential of simply quantifying and accounting for changing user intent for extending existing systems. Note also that the discussions are easily extendable to other contexts with sequential interactions, such as music or news recommendations, in line with the central goal of adaptability.

While we do propose a related evaluation metric through this new lens of relevance, by comparing the counter-factual *DispAdd* from the recommended next item to the ground truth, this work differs in spirit from the broader study of diversification in recommender systems ([8, 9, 10]). Rather than making a prescriptive claim on a proposed metric and proposing an algorithm that improves performance with respect to the metric, e.g. novelty, our focus is more descriptive and aims at dynamic user profiling.

The paper is organized as follows. We discuss comparisons with existing literature in Section 2. In Section 3, we present the detailed definitions and properties of our proposed concepts, illustrated through analysis of data collected from Etsy. Then in Section 5, we show proof-of-concept experimental results using three e-commerce datasets for next-item recommendation and purchase prediction. Finally, we conclude with Section 6.

To summarize, our main contributions are as follows:

- We quantify *affinity to variation* in sequential recommendation to capture dynamic user exploration preferences, offering intuitive ingredient for adaptive approaches.
- We establish soundness and usability of the proposed concepts through analysis both mathematically and based on real-world e-commerce data.
- We demonstrate flexibility and potential by enhancing baseline algorithms with proposed concepts for two separate use cases, showing improved performance in three datasets.

2. Related Work

The dynamic nature of recommender systems has long been recognized. Sequential recommender models that make personalized next-item predictions based on interaction histories rose rapidly in popularity, and form the context of the discussion here. Well-known models include GRU4REC [11], Caser [12], SASRec [5], NextIt[13], BERT4Rec [6], S3-Rec [7]; see surveys [14, 15, 16].

The discussion of exploration preference in this work is based on the fundamental recognition that similarity to past interactions alone does not warrant user satisfaction. Empirical evidence supporting this notion [17] in part spurred development of the field of diversification in recommender systems, concerned with evaluation metrics beyond accuracy. Most common ones include diversity [1, 18], novelty [19, 20], and serendipity/ unexpectedness [21, 22, 23]; we defer reader to surveys for details [24, 25, 26, 27, 28]. This work share similar motivations but focuses more on describing user behavior, compared to the more prescriptive nature of studies that centers around promoting alternative objectives that are important for recommender quality.

Nevertheless, many recent works take a personalized approach that bear relevance to this work or could be inspiration for future directions. Some approaches balance the tradeoff between accuracy and alternative objectives based on user behavior [8, 29, 30, 31, 32]. User demonstrated preference is sometimes considered explicitly, in a similar spirit to our work. Kapoor et al. [33] propose the concept of novelty preference which changes based on user and session, and is predicted with the help of behavioral psychology insights. Li et al. [34] use multi-cluster modeling of user interests in the latent space to personalize unexpectedness using self-attention. Mehrotra et al. [35] log user behavioral response to divergent recommendations to predict future receptivity. Qian et al. [36] disentangle user intent into popularity conformity and personal preference to introduce intrinsic novelty for long-tail items.

Another aspect of our proposed quantification involves the dependency on context that changes over time. While it is mostly implicit here through user behavior, many recent works explicitly model these temporal context and effects [37, 38, 39, 40, 41, 42, 43]. Of particular interest is the concept of *temporal diversity* proposed by Lathia et al. [44] that similarly explores the dissimilarity between sets of top-N recommendations temporally, though they concern mainly with the same items repeating over time.

Lastly, we highlight some studies on the particular effect of user intent on preference drift [45, 46, 47]. Intuitively, the exploration preference of an user is linked to how strongly they intend to purchase, and indeed we use purchase prediction as one of the example use cases.

However, the central focus of this work is not on determining what governs user intent but rather on representing and adapting to exhibited behavior.

3. Affinity to Variation

We now provide formal definitions of the concepts described thus far, offering design justifications and useful properties.

Recall again that we are interested in a given user’s *affinity to variation* given their recent interactions, based on how concentrated the items when embedded on a vector space. We assume access to the sequential interaction history of each user, where each interaction consists of an item and an action. Furthermore, we assume access to a reasonable embedding space, as detailed later.

Dispersion For a fixed user, consider a sequence of interactions, $\{(x_1, a_1), (x_2, a_2), \dots, (x_T, a_T)\}$, where x_t is the item the user interested with at time t and a_t is the action taken. Denote this sequence of items in the embedding space as $\{X_1, X_2, \dots, X_T\}$.

Let L be the size of the look-back window, where we consider the last L interactions to be a snapshot of a user’s current attention space. To assess the concentration level of this attention space, we measure the spread of these items on the embedding space using high-dimensional analogue of variance, known in some contexts as the squared standard distance deviation.

We define the *Dispersion* of the L -step sub-sequence ending at time step t as

$$\sigma^2(t; L) := \frac{1}{L} \sum_{i=t-L}^t \|X_i - \mu(t; L)\|^2,$$

where $\mu(t; L) = \frac{1}{L}(\sum_{j=t-L}^t X_j)$ and $\|\cdot\|$ is the Euclidean norm. In other words, *Dispersion* measures the average squared Euclidean distance of the points in the set from the mean of the set. We omit L when it is clear from the context.

This is meant to be a simple and convenient definition, aiming towards our ambitions of usability in practice and adaptability for future research. The simplicity of the definition, along with its resemblance to variance, helps with intuition and interpretation. Moreover, similar to the standard definition of variance, *Dispersion* can be updated incrementally as the sequence rolls forward, without having to store every embedding in the sequence and thereby reducing memory demands. Having a simple construct also help the concepts to be well-behaved, and thus conducive to future adaptations. Specifically, the standard components of *Dispersion* lead to well-understood properties that allow easy incorporation into deep neural networks, while the problem-agnostic range of *DispAdd* makes it viable for settings such as weighting or hypothesis testing.

DispAdd With this basic quantification of how spread out the user’s interactions are over a fixed window, we examine the dynamics of how the spread changes with each new interaction.

With each subsequent item that the user interacts with, the window of L most recent items is shifted forward to include this new item. Because the window size is fixed, how the relative distance from the set to the new item compares to that to the least recent item is reflected in the change in *Dispersion*.

While direction of change is the most important for detecting trends, we would also want to understand the strength of the trend through magnitude of change. However, note that the scale of *Dispersion* depends in general on the distribution of the items on the embedding space, and specifically on which regions the recent set of interaction items are mapped to, so the difference in *Dispersion* could vary greatly in scale. This implies concerns in interpretability not only from simple difference but also from percentage change, since the resulting *Dispersion* could be orders of magnitude larger than the preceding one, especially for embedding space that is normalized.

To capture this change while being scale invariant, we define *DispAdd* as the relative difference resulting from shifting the window, comparing the difference to the average:

$$DispAdd(t) := \frac{\sigma^2(t) - \sigma^2(t-1)}{\frac{1}{2}(\sigma^2(t) + \sigma^2(t-1))}. \quad (1)$$

A positive value of *DispAdd* indicates an increase in *Dispersion*, reflective of increased affinity to variation; the reverse is true of a negative value.

Note that the risk of division by zero is reduced, but not eliminated, in which case the definition becomes ill-posed. This only happens in the special case that $\sigma^2(t) = \sigma^2(t-1) = 0$; there defining *DispAdd* to be 0 seems reasonable. However, for the purpose of our discussions, we will in fact disregard this case. A *Dispersion* value of 0 implies that all items interacted in the look-back window are at the same point on the embedding space, or identical from the view of the vector representation. While this could be possible, it could also indicate an embedding space without enough representation power, unintended system behavior or even malicious bot activity. It is more informative for raising concerns than our purpose of understanding user behavior patterns.

Putting aside the above technicality, the range of *DispAdd* is bounded independent of the embedding space, as stated below.

Lemma 3.1. *DispAdd has value bounded between -2 and 2.*

The proof is in the appendix. This is a theoretical range to certify the bounded behavior of the definition, while in practice the real range could be much narrower as shown in the experiments, depending on the characteristics of the embedding space.

ADAD We treat *DispAdd* as a dimension of user intent, and evaluate a sequential recommender system in its ability to correctly capture it. We argue that recommended items that result in similar *DispAdd* as the true next-item better match the user’s intended behavior, reflecting an alternative aspect of relevance.

The time indices are less important for evaluating a recommender system, as recommender outputs to predefined test trajectories are compared to the ground truths. Thus we instead index trajectories to simplify notation. Specifically, let $DispAdd_j$ denote the relative change to

Dispersion from the true next item following trajectory j . Then replace the ground truth by the next-item predicted by the recommender, and denote the hypothetical *Dispersion* as $\bar{\sigma}_j^2$, and hypothetical relative change using $\bar{\sigma}_j^2$ as $\overline{DispAdd}_j$. The absolute difference between the two values represents the error in capturing the user’s true change in affinity to variation.

To evaluate on the recommender level, we average the above difference over all trajectories in the test set \mathbb{T} , and call it Average DispAdd Deviation (*ADAD*):

$$ADAD := \frac{1}{|\mathbb{T}|} \sum_{j \in \mathbb{T}} (|\overline{DispAdd}_j - DispAdd_j|).$$

In the ideal case that the correct *DispAdd* is perfectly captured and the error is 0, (such as is the case if the predicted item is the true next item), then $\bar{\sigma}^2(j) = \sigma^2(j)$.

In fact, we show this two way relation is true in general:

Lemma 3.2. $|\overline{DispAdd}_j - DispAdd_j| = 0$ iff $\bar{\sigma}^2(j) = \sigma^2(j)$.

Consequently, we will explicitly model $\bar{\sigma}^2(j)$ and $\sigma^2(j)$ instead of the respective relative differences in the next section, for simplicity and numerical stability.

Embeddings The above definitions rest on the assumption that there is an embedding that can sufficiently capture the relation among items, specifically that pairwise distances between item embeddings have reasonable correspondence to the similarities between the items. How to encode rich item information appropriately into embedding vectors is an actively research area in itself, with no universally agreed upon best practice. On the other hand, many recommender systems train their own embedding mappings in the process of learning item relevance, where the performance of the system lends credibility to the quality of the embeddings. Our proposed definitions work with embedding either trained separately or from a base recommender engine.

4. Analysis of platform data

We now explore our definitions of *Dispersion* and *DispAdd* on real-world data. We draw samples from Etsy user interaction data, and calculate embedding vectors using the Universal Sentence Encoder [4]. We present results from moderately active users, specifically those who made between 2 and 6 purchases in the past year, to avoid any potential idiosyncrasies of the extremes. Results from highly active users in fact show similar patterns. We sample a percentage from the candidate pool, roughly 3,000 users, and collect their interactions from the first six months of 2021. There are 4 types of actions: view, favorite, cart, and buy. Users with less than 20 interactions and 3 purchases are filtered, and consecutive interactions within 4 hours of each other that are identical are removed. The resulting dataset contains 799 users, with an average of 164 interactions per user.

First, we illustrate the definitions with example trajectories to provide a concrete understanding, and to check for consistency with our intuition. We calculate *Dispersion* for a window of six items, displayed on the top row in each figure, then calculate the *DispAdd* for each of the four items that follow, shown on the bottom row.

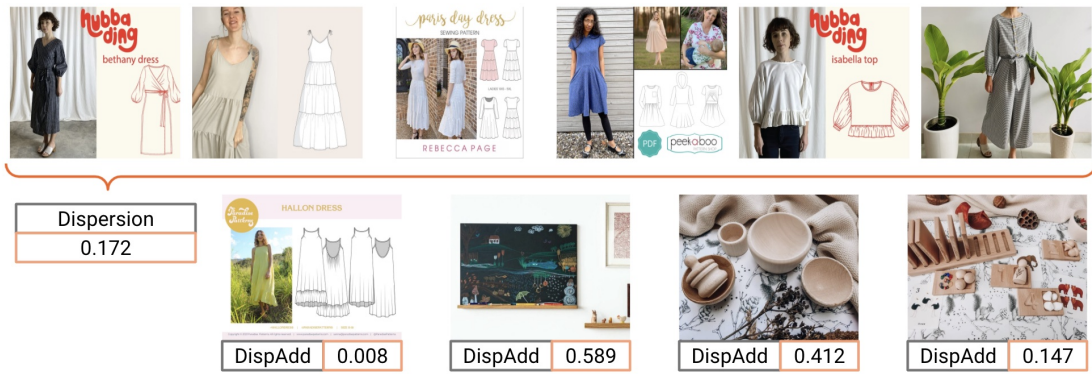


Figure 2: Example trajectory 1 with *Dispersion* over a 6-item window and subsequent *DispAdd* displayed.



Figure 3: Example trajectory 2 with *Dispersion* over a 6-item window and subsequent *DispAdd* displayed.

In Figure 2, we can see that the user starts by interacting with sewing patterns of similar styles; this is corroborated by a small *Dispersion* value of 0.172. The first item that follows is again similar, confirmed by a near-zero *DispAdd* that reflects little change. However, each of the next two items belongs to a different category, showing an expanding affinity to variation; indeed, we see large positive *DispAdd* values for both steps. The last item is less of a departure from the item prior, showing a slow in the exploration, which the small positive *DispAdd* value again confirms.

The opposite pattern can be seen in Figure 3. Here the user begins with a diverse selection including various kitchen items and notebooks, captured by a relatively large *Dispersion* value of 0.507. The spread then shrinks as the user focuses in on notebooks of a particular style, during which process we see increasingly negative *DispAdd* values. Then, when the user interacts with an item from a different category, there is a slight expansion in the affinity to variation again, reflected by a small positive *DispAdd*.

Now, we attempt to verify intuitive hypotheses about user behavior patterns using *DispAdd*. Specifically, it seems reasonable to expect that before a user makes a purchase, they have a

higher tendency to look at items that are highly similar to each other to compare. After the purchase is complete, on the other hand, we might expect them to be more interested in items that are meaningfully different, as their previous need has been met. To extract supporting evidence, we look at each time a purchase is made, and take the last five interacted items prior to it as well as the five following it. We compute the *DispAdd* resulting from each interaction, and average them according to the number of interactions from the purchase (negative corresponds to before purchase and positive corresponds to after); we plot the results in Figure 4 for four different window sizes. We see that there is a clear distinction between items interacted with prior to purchase, which produced negative *DispAdd* on average, and items interacted with after a purchase, whose resulting *DispAdd* increased noticeably. The trends are consistent for the window sizes presented. Note that smaller window sizes result in less smoothness and vice versa, so we choose sizes which are sufficiently large but not so much to cause overly muted dynamics.

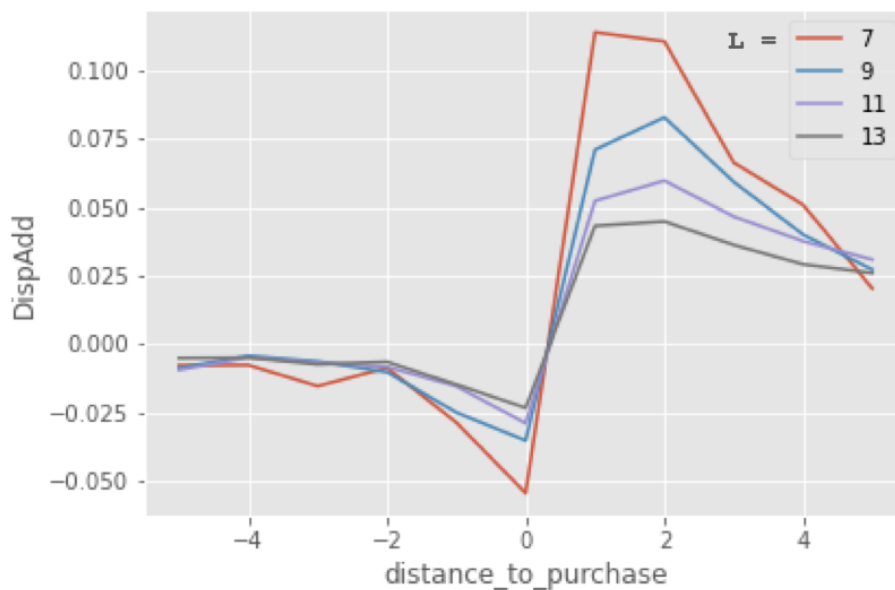


Figure 4: Average *DispAdd* for items interacted with immediately before and after a purchase, plotted for different window sizes.

We repeat a similar analysis, but divide the behaviour relative to the start of a new session, which is an interaction that is at least 4 hours apart from the previous one. We again take five items from the end of the previous session and five items from the start of the new session, and average the *DispAdd* from each interaction. We observe a similar trend, though of a reduced scale. In Figure 5, we again plot the average *DispAdd* against the distance, defined similar to above. The plots reflect a tendency for users to look over items very similar to recent interactions right before the end of a session, and interact with items further removed from the average of recent interactions at the start of a new session. Compared to the behavior after a purchase, this exploration seems to end much sooner on average, indicated by a quick drop of

the average *DispAdd* towards 0. This aligns with user interviews conducted by the platform, where users report themselves typically returning via the home page, where they may engage briefly with prompts for exploratory contents, then go back to where they left off.

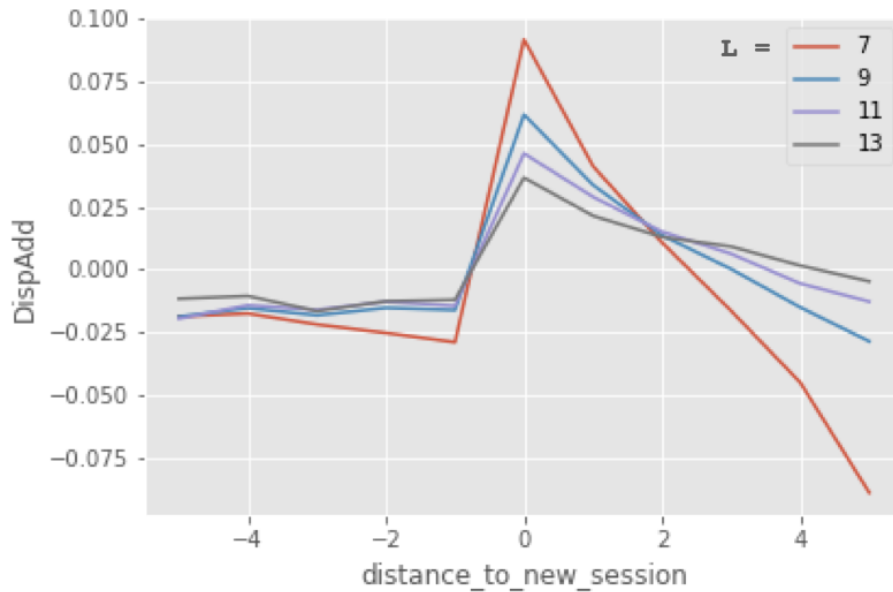


Figure 5: Average *DispAdd* for items interacted with near the end of a session and at the start of a new one, plotted for different window sizes.

We now examine patterns on the taxonomy level. Etsy categorizes products with 15 top-level taxonomies; for each interaction, we average its *DispAdd* into the taxonomy the item belongs to, and present the resulting bar plot on the bottom of Figure 6. To contrast with alternative quantification without the new concepts, we plot on the top of Figure 6 the average number of consecutive interactions in the same taxonomy (blue solid line) and the total number of purchases of this taxonomy (red dashed line).

When a taxonomy has negative average *DispAdd*, it suggests that users overall spend more time in narrowing behavior when interacting with items in the category. In other words, the plot shows listings under the categories of Jewelry and Wedding, for example, are more likely to be visited following interactions with similar items. This largely aligns with our intuitive understanding of user behaviors. For jewelry, one typically moves through many iterations of products that have increasing number of desirable features, easily resulting in long trajectories of increasingly similar items. This is consistent with the alternative statistics, with largest number of consecutive interactions of all taxonomies. For wedding supplies, it is reasonable to imagine that most users engaging with the category have persistent interest over a prolonged period of planning, and have generally static style preferences. The prolonged planning is confirmed by the low total purchases, and the short average consecutive interactions may be due to interactions with other categories in-between.

On the other hand, when a taxonomy has positive average *DispAdd*, it reflects that users visit the category more frequently as part of exploratory behavior rather than focused engagements.

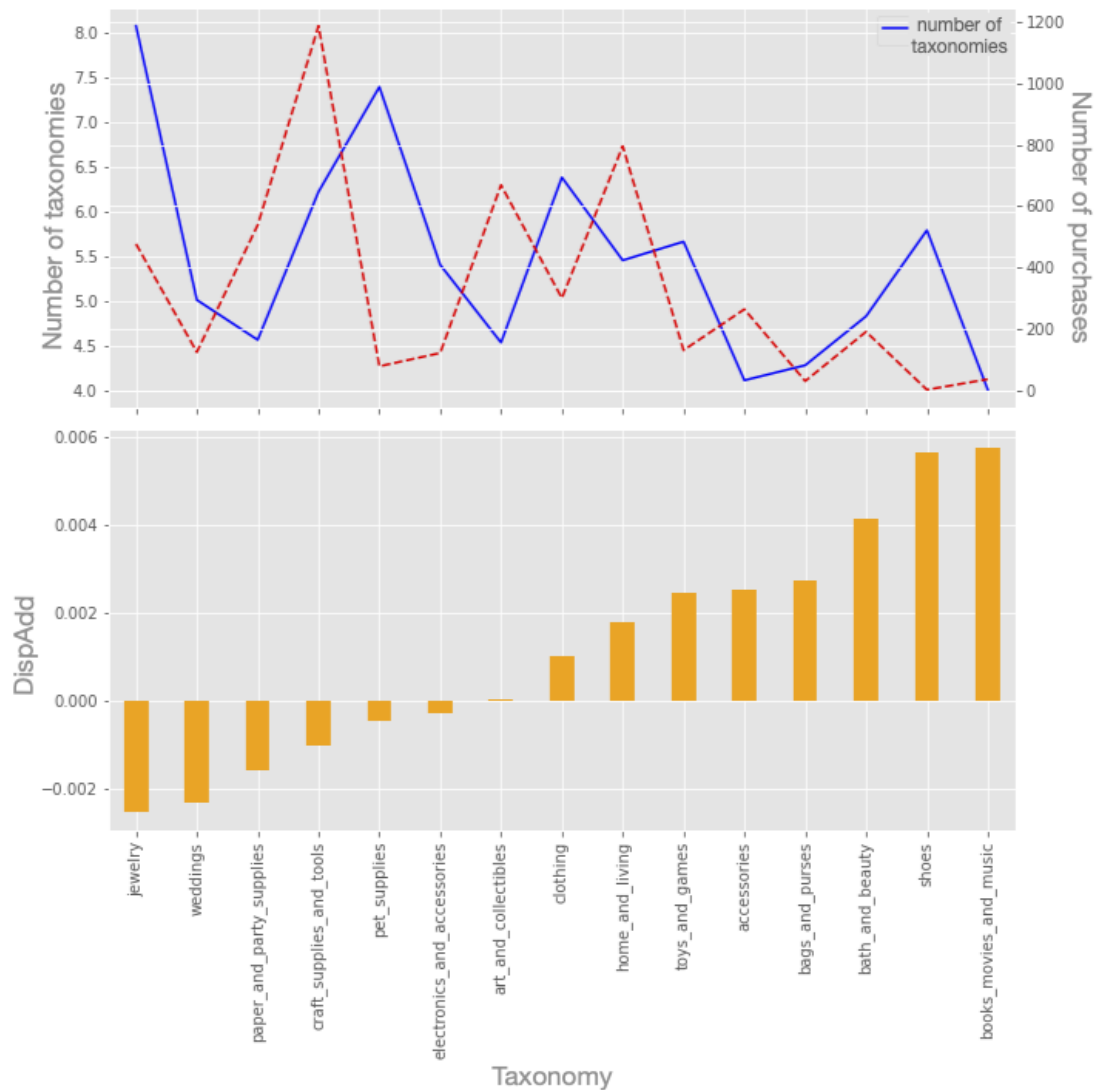


Figure 6: Etsy dataset statistics broken down by the top level taxonomies. (TOP): average consecutive interactions (blue solid line) and total purchases (red dashed line) in the taxonomy; (BOTTOM): average *DispAdd* of items belonging to the taxonomy.

This usually mean less purchasing, as confirmed by the overall lower total purchases for taxonomies with more positive *DispAdd*. There is one category, Home and Living, that stands out as having relative a high purchase number despite the positive average *DispAdd*. One potential interpretation might be that highly similar items are less common in the category, or that many items serve some specific functionalities, so there is less incentive to re-visit the same or a similar item once a user is satisfied with a discovery.

In both cases, we see that alternative statistics are able to help interpret patterns uncovered

Table 1
Dataset statistics (after pre-processing)

Dataset	# users	# items	avg actions / user	# actions
Retail-Rocket	35,423	29,437	11.99	495,400
Tmall	8,133	497,022	694.41	5,663,878
Etsy	23,441	2,008,696	163.13	3,870,776

by affinity to variation, but not reflect these patterns themselves. Our proposed concepts, in reflecting this dynamic form of intent otherwise hard to capture, enables a novel angle for gaining insights into user behavior.

5. Numerical Evidence

The ability to quantify a user’s affinity to variation opens up a range of application opportunities. For proof-of-concept, we incorporate the proposed concepts in two common applications: next-item recommendation and purchase prediction.

5.1. Datasets

We investigate the performance on three different e-commerce datasets, two publicly available and one proprietary. In general we remove users and actions with too few interactions to ensure a base level of density. Moreover, since we aim to capture changing user behavior patterns, we only preserve users with a minimum number of purchases. The data cleaning procedures are detailed below, and the resulting summary statistics are in Table 1.

Retail-Rocket: An e-commerce dataset publicly available on Kaggle ¹ that contains 4.5 months of user interaction data. Actions include click, add-to-cart, and transaction. We set minimum number of interactions to 5, but no minimum purchase due to data size.

Tmall: A publicly available dataset provided by Alibaba ² that contains around 10,000 user records and 12 million actions of user activities on Taobao app from Nov 18, 2014 to Dec 18, 2014. Each action is one of click, add-to-favorite, add-to-cart and purchase. We set interaction minimum to 5, and purchase minimum to 2.

Etsy: A proprietary dataset collected from user activity logs between Aug 1, 2021 to Aug 31, 2021. A random sample is drawn from the most active users on the platform, containing the user id, item id, and interaction type, which include view, cart, favorite and purchase. Minimum number of interactions is set to 20, and minimum purchase is set to 3.

For each user trajectory, the second to last item is used for validation and the last for testing. All remaining interactions with at least L steps of history are used for training. Hyperparameters are tuned on the validation sets, then results on the test sets are reported. For more efficient evaluation, for each target item we sample 99 negative examples to create the candidate set,

¹<https://www.kaggle.com/retailrocket/e-commerce-dataset/home>

²<https://tianchi.aliyun.com/dataset/dataDetail?dataId=46>

and evaluate on a sample of 10,000 users when applicable. When an identical action is repeated, only the first is preserved to reduce noise.

5.2. Embeddings and Base Engine

We opt to use embeddings trained with a base recommendation engine to demonstrate the flexibility of our construct. This is particularly convenient for recommendation tasks such as our first example, since *Dispersion* can simply be incorporated in an additional module on top of the base model.

For clarity and consistency, we choose to use one base engine for both experiments, varying the size and sparsity of the datasets instead. This of course raises the concern of generalizability. Our ideal candidate is a base engine with verified performance to ensure the informativeness of our defined concepts, but not overly optimized or specialized, in order that the outcomes of the experiments can be reflective of similar engines. We examined the performance of two popular baseline methods, SASRec ([5]) and GRU4Rec ([11]), both applicable to general settings with no requirements on data density or context information. We found that the characteristics of *Dispersion* differed more among the datasets of varying size and sparsity using the same embedding mapping than between the embedding mappings, supporting the choice to prioritize varying the datasets in the experiments.

SASRec is what we present here since it is a representative sequential recommendation approach. It is used as baseline in numerous studies with consistent performance (e.g., [13, 41]) and extended upon in later works (e.g., [6, 46]). In fact, the attention mechanism it uses from the Transformer structure ([48]) is actively studied and extended upon both in and outside of recommender systems ([49, 50]), proving to be a adaptable method with the potential to persist as a component in a subclass of future state-of-the-art methods. This in turn hopefully lends more relevance to the demonstrated performance boosts from incorporating our proposed concepts.

5.3. Application 1: Next-item Recommendation

As affinity to variation is defined in the context of sequential interactions, next-item recommendation is the most natural application.

From the exploratory analysis, we see that the pattern of a user’s engagement is not static across time. For example, one might engage with items similar to the previous one for some period of time while opting for entirely different categories during others. While many sophisticated sequential models have been proposed to account for various short-term and long-term trends (e.g. [51, 52]), *DispAdd* provides direct information on the change in the user’s affinity to variation. Thus, we aim to augment existing sequential recommender systems by also considering the *Dispersion* signal.

Specifically, in order to choose the most relevant item to recommend next, we predict *Dispersion* at the next time step and assign a score to each candidate item based on how closely their inclusion would match this predicted level of *Dispersion*. In particular, to recommend a next-item at time T , we use item embeddings from the base model to compute $\sigma^2(t)$ for $t \leq T$, and predict $\hat{\sigma}^2(T + 1)$. Then, a relevance score for each of the candidate items

is computed based on the closeness to generating $\hat{\sigma}^2(T + 1)$. Here we do so by computing $\bar{\sigma}^2(T + 1)$ induced by each candidate item and then using the negative of the L1 distance to $\hat{\sigma}^2(T + 1)$ as logits. We perform the prediction with a simple two-layer network, as illustrated in Figure 7 (LEFT), but any prediction model (and qualifying base engine) can be swapped into the general framework, indicated by the dashed box(es) in Figure 7 (RIGHT).

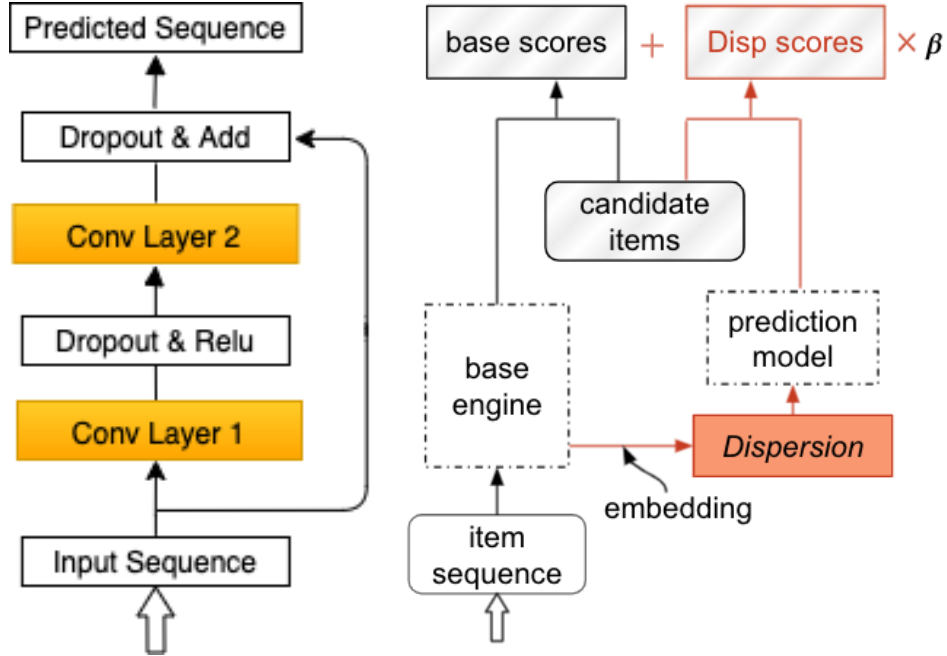


Figure 7: (LEFT) Structure of a simple two-layer feedforward network; "Add" refers to residual connection. (RIGHT) General framework for incorporating *Dispersion* in next-item recommendation.

In general, a base engine qualifies if it (1) is a sequential recommendation model, (2) uses embedding vector representations for items, and (3) generates recommended next items from a set of candidates by computing their relevance scores. These are common to some of the most popular models for sequential recommendations, ensuring that the structure is generalizable and expandable.

Note that a possible alternative is to implement a secondary prediction algorithm for the desirable item embedding that would generate *Dispersion* of value $\hat{\sigma}^2(T + 1)$, then compute the relevance with inner product, which might be more efficient for large systems.

With the new relevance scores based on affinity to variation, we can obtain a different ranking of the candidate items by adding them to the original scores produced by the base engine. To ensure the relative magnitude of the two is appropriate, we scale the new scores to match the mean of the original scores, then weigh by a changeable hyperparameter β . We rank the candidate items using this combined score.

We adopt the hyperparameters and initialization strategies suggested by the author for SASRec. We use learning rate 0.001, latent dimension 50, l2 regularization parameter 1e-6, and maximum sequence length roughly proportional to the average action length: Retail-Rocket at

Table 2
Next-Item Recommendation

Dataset	Method	HR@1	HR@5	NDCG@5	ADAD (\downarrow)
R-R	base	0.6838	0.8393	0.7700	0.2052
	combined	0.6928	0.8397	0.7737	0.2028
Tmall	base	0.7519	0.8436	0.8033	0.1018
	combined	0.7584	0.8437	0.8057	0.1010
Etsy	base	0.4387	0.5089	0.4785	0.2025
	combined	0.4626	0.5088	0.4877	0.1996

30, Tmall at 300 and Etsy at 100. We tuned the hyperparameters for the *Dispersion* prediction module on the validation set. We set the learning rate to 0.001 after performing a grid search from $\{0.1, 0.01, 0.001, 0.0001\}$, and the weight in the combined score to 0.8. We set the window/kernel size to 6/4, 8/5, 7/3 for Retail-Rocket, Tmall, and Etsy respectively, from $\{6,7,8,9,10\}$ for window and $\{1,2,3,4,5,6\}$ for kernel. The reported results are based on outputs for the test set using the combined score between the tuned *Dispersion* module and SASRec. However, note that improvements are observed consistently for other hyperparameter choices as well, without discernible trends.

To evaluate the performance, we compare the ranking resulting from the combined scores to the ranking using only the base scores. Table 2 contains the results, reported as the average over 5 repeated runs; results that are statistically better, with p-value less than 0.05, are in bold. Retail-Rocket is listed as "R-R".

For accuracy, we report popular measures HitRate@5, NDCG@5 and HitRate@1. Across all datasets, we see that the combined score produced rankings with higher NDCG@5 and HitRate@1 compared to the base model alone, showing a boost from incorporating *Dispersion* information even in this simplest form. We also report the *ADAD* values, which are lower for rankings using the combined score. Since we argue a recommendation that better matches the *DispAdd* of the true next-item is more relevant to the user, it suggests additional reduction in user dissatisfaction not captured by the improved accuracy.

Note that for comparison, we also compute *ADAD* for a randomly selected item to follow each trajectory, which is generally larger but never exceeds 0.6. It seems that the embedding space generally produces small differences in *DispAdd* for the item sets. Furthermore, this indicates the results of the base recommender on average matches the users' demonstrated *DispAdd* much better than random, which is as we expect with the accuracy performance.

5.4. Application 2: Predicting Purchases

A second example application we explore is the connection between affinity to variation dynamics and purchase behaviors. Analysis presented before shows an overall narrowing trend in interest sets prior to purchases, so a natural question is whether modeling *Dispersion* might help forecast whether a purchase will occur.

Answering this question is meaningful for e-commerce platforms, as it could help understand user intents and better cater to real time needs. This would be especially the case if a reliable purchase forecast could be obtained a period of time in advance (e.g. at the start of a session).

Table 3
Purchase Prediction

	Retail-Rocket		Tmall		Etsy	
	AUC	PR-A	AUC	PR-A	AUC	PR-A
action	0.9232	0.4010	0.6472	0.0781	0.7270	0.4279
+ disp	0.9295	0.4213	0.6560	0.0866	0.7289	0.4304

As a first step, we focus on a simple sequential setting, and aim to predict whether the next interaction will be a purchase using information up to the current time.

The most basic way to attempt this is to perform a time series forecast using actions up to the current time as features. This is the approach we compare against as the baseline. To investigate the value from information of affinity to variation, we simply include *Dispersion* as an additional feature at each time step. In other words, now there are two covariates for the timeseries forecast. Note that *Dispersion* is used as the feature instead of *DispAdd* because it preserves the magnitude information, while the key signal from *DispAdd*, the sign, can be derived from the first difference.

We use the same two-layer network as shown in Figure 7, changing only the input and output dimensions in addition to tuneable hyperparameters. For fairness of comparison, we tune the hyperparameters based on the baseline model using only action signals, then keep the same model configuration when adding the *Dispersion* signal. For each, the model with the best performance on the validation set during training is selected and the performance on the test set reported. Because the number of purchases is much lower than the number of non-purchases, we report not only the AUC (area-under-the-curve) value of the Receiver Operating Characteristic (ROC) curve, which shows the general diagnostic ability of the binary classifier, but also the AUC value for the Precision-Recall (PR) curve, which shows the trade-off between precision and recall of only the positive class.

The results for the three datasets are shown in Table 3, with ROC-AUC reported as "AUC" and Precision-Recall AUC reported as "PR-A". The learning rate used for Retail-Rocket and Etsy is 0.01, and for Tmall is 0.1, after performing a grid search from {0.1, 0.01, 0.001, 0.0001}. The same window and kernel sizes are used, and the number of hidden units in the feedforward network is set to 16 for Retail-Rocket and 32 for the remaining two datasets, after searching {2, 4, 8, 16, 32}. We observe that the prediction accuracy for the purchase action is higher for all three datasets. While the improvements seem marginal for ROC-AUCs, the percentage improvements for PR-AUCs are much higher, reflecting stronger performance in precision and recall for positive instances of purchase actions despite their rare occurrence. Note that the predictive ability for Tmall is poor overall due to an extremely small number of purchase signals in the test set, only 1.1% at 91 purchases.

5.5. Future Directions

Our research highlights the importance of understanding users' exploration preference dynamically, and establishes a foundation for new approaches in adaptive systems. One could envision an ensemble recommender system with specialized context dependent modules, modulated

based on the user's current affinity to variation. This could help open up possibilities on more creative approaches to recommendations, such as a style based exploration module. If a large increase in affinity to variation is anticipated, the module is activated to offer recommendations that match recent interactions in style but not necessarily category, providing more opportunities for serendipitous discoveries. As a further potential for interplay with diversification work, our derived metric could help ground the appropriate balance between accuracy and diversification objectives. By demonstrating user demonstrated preference, it takes away the burden of arguing for diversification to be added at the "cost" of accuracy, and introduces an additional dimension to the inherently multi-faceted concept of relevance.

Other future directions include alternative constructs for user's exploration preference, such as using determinant of embedding matrix, or more dynamic treatment of the lookback window size. It may also be worth examining the confounding effect of existing recommender policies on observed user behavior regarding exploration preferences, which could yield behavioral insights and present opportunity for further studies on how to remove such bias.

6. Conclusion

In this paper, we introduced novel concepts to quantify and analyze user exploration preference in the context of sequential interactions. We propose affinity to variation measure at each time step, *DispAdd* to measure the relative change, and a derived metric to evaluate whether results from a recommender system matches the true *DispAdd*. These definitions with their mathematical properties establish a foundation for building adaptive recommendation systems that account for the dynamic nature of user preferences, demonstrated through promising results in proof-of-concepts experiments. The insights and methodologies presented in this paper open up simple and implementable avenues to enhance recommendation algorithms for more relevant item suggestions to users. With potential to connect ongoing research, we believe that our approach opens exciting new directions for progress in more adaptive and effective recommendation systems in the future.

References

- [1] K. Bradley, B. Smyth, Improving recommendation diversity, in: Proceedings of the twelfth Irish conference on artificial intelligence and cognitive science, Maynooth, Ireland, volume 85, Citeseer, 2001, pp. 141–152.
- [2] S. M. McNee, J. Riedl, J. A. Konstan, Being accurate is not enough: how accuracy metrics have hurt recommender systems, in: CHI'06 extended abstracts on Human factors in computing systems, 2006, pp. 1097–1101.
- [3] C. Demangeot, A. J. Broderick, Exploration and its manifestations in the context of online shopping, *Journal of Marketing Management* 26 (2010) 1256–1278.
- [4] D. Cer, Y. Yang, S.-y. Kong, N. Hua, N. Limtiaco, R. S. John, N. Constant, M. Guajardo-Céspedes, S. Yuan, C. Tar, et al., Universal sentence encoder, *arXiv preprint arXiv:1803.11175* (2018).

- [5] W.-C. Kang, J. McAuley, Self-attentive sequential recommendation, in: 2018 IEEE International Conference on Data Mining (ICDM), IEEE, 2018, pp. 197–206.
- [6] F. Sun, J. Liu, J. Wu, C. Pei, X. Lin, W. Ou, P. Jiang, Bert4rec: Sequential recommendation with bidirectional encoder representations from transformer, in: Proceedings of the 28th ACM international conference on information and knowledge management, 2019, pp. 1441–1450.
- [7] K. Zhou, H. Wang, W. X. Zhao, Y. Zhu, S. Wang, F. Zhang, Z. Wang, J.-R. Wen, S3-rec: Self-supervised learning for sequential recommendation with mutual information maximization, in: Proceedings of the 29th ACM international conference on information & knowledge management, 2020, pp. 1893–1902.
- [8] M. Zhang, N. Hurley, Novel item recommendation by user profile partitioning, in: 2009 IEEE/WIC/ACM International Joint Conference on Web Intelligence and Intelligent Agent Technology, volume 1, IEEE, 2009, pp. 508–515.
- [9] R. Boim, T. Milo, S. Novgorodov, Diversification and refinement in collaborative filtering recommender, in: Proceedings of the 20th ACM international conference on Information and knowledge management, 2011, pp. 739–744.
- [10] J. Parapar, F. Radlinski, Diverse user preference elicitation with multi-armed bandits, in: Proceedings of the 14th ACM International Conference on Web Search and Data Mining, 2021, pp. 130–138.
- [11] B. Hidasi, A. Karatzoglou, L. Baltrunas, D. Tikk, Session-based recommendations with recurrent neural networks, arXiv preprint arXiv:1511.06939 (2015).
- [12] J. Tang, K. Wang, Personalized top-n sequential recommendation via convolutional sequence embedding, in: Proceedings of the eleventh ACM international conference on web search and data mining, 2018, pp. 565–573.
- [13] F. Yuan, A. Karatzoglou, I. Arapakis, J. M. Jose, X. He, A simple convolutional generative network for next item recommendation, in: Proceedings of the Twelfth ACM International Conference on Web Search and Data Mining, 2019, pp. 582–590.
- [14] C. Rana, S. K. Jain, A study of the dynamic features of recommender systems, *Artificial Intelligence Review* 43 (2015) 141–153.
- [15] M. Quadrana, P. Cremonesi, D. Jannach, Sequence-aware recommender systems, *ACM Computing Surveys (CSUR)* 51 (2018) 1–36.
- [16] I. Rabiou, N. Salim, A. Da’u, A. Osman, Recommender system based on temporal models: a systematic review, *Applied Sciences* 10 (2020) 2204.
- [17] C.-N. Ziegler, S. M. McNee, J. A. Konstan, G. Lausen, Improving recommendation lists through topic diversification, in: Proceedings of the 14th international conference on World Wide Web, 2005, pp. 22–32.
- [18] S. Vargas, L. Baltrunas, A. Karatzoglou, P. Castells, Coverage, redundancy and size-awareness in genre diversity for recommender systems, in: Proceedings of the 8th ACM Conference on Recommender systems, 2014, pp. 209–216.
- [19] N. Hurley, M. Zhang, Novelty and diversity in top-n recommendation—analysis and evaluation, *ACM Transactions on Internet Technology (TOIT)* 10 (2011) 1–30.
- [20] S. Vargas, P. Castells, Rank and relevance in novelty and diversity metrics for recommender systems, in: Proceedings of the fifth ACM conference on Recommender systems, 2011, pp. 109–116.

- [21] Y. C. Zhang, D. Ó. Séaghdha, D. Quercia, T. Jambor, Auralist: introducing serendipity into music recommendation, in: Proceedings of the fifth ACM international conference on Web search and data mining, 2012, pp. 13–22.
- [22] P. Adamopoulos, A. Tuzhilin, On unexpectedness in recommender systems: Or how to better expect the unexpected, *ACM Transactions on Intelligent Systems and Technology (TIST)* 5 (2014) 1–32.
- [23] D. Kotkov, S. Wang, J. Veijalainen, A survey of serendipity in recommender systems, *Knowledge-Based Systems* 111 (2016) 180–192.
- [24] S. Gollapudi, A. Sharma, An axiomatic approach for result diversification, in: Proceedings of the 18th international conference on World wide web, 2009, pp. 381–390.
- [25] M. Ge, C. Delgado-Battenfeld, D. Jannach, Beyond accuracy: evaluating recommender systems by coverage and serendipity, in: Proceedings of the fourth ACM conference on Recommender systems, 2010, pp. 257–260.
- [26] M. Kaminskis, D. Bridge, Diversity, serendipity, novelty, and coverage: a survey and empirical analysis of beyond-accuracy objectives in recommender systems, *ACM Transactions on Interactive Intelligent Systems (TiiS)* 7 (2016) 1–42.
- [27] T. Silveira, M. Zhang, X. Lin, Y. Liu, S. Ma, How good your recommender system is? a survey on evaluations in recommendation, *International Journal of Machine Learning and Cybernetics* 10 (2019) 813–831.
- [28] Q. Wu, Y. Liu, C. Miao, Y. Zhao, L. Guan, H. Tang, Recent advances in diversified recommendation, *arXiv preprint arXiv:1905.06589* (2019).
- [29] N. Tintarev, M. Dennis, J. Masthoff, Adapting recommendation diversity to openness to experience: a study of human behaviour, in: *International Conference on User Modeling, Adaptation, and Personalization*, Springer, 2013, pp. 190–202.
- [30] C. H. Teo, H. Nassif, D. Hill, S. Srinivasan, M. Goodman, V. Mohan, S. Vishwanathan, Adaptive, personalized diversity for visual discovery, in: Proceedings of the 10th ACM conference on recommender systems, 2016, pp. 35–38.
- [31] Y. Song, N. Sahoo, E. Ofek, When and how to diversify—a multicategory utility model for personalized content recommendation, *Management Science* 65 (2019) 3737–3757.
- [32] R. S. Fortes, D. X. de Sousa, D. G. Coelho, A. M. Lacerda, M. A. Gonçalves, Individualized extreme dominance (inded): A new preference-based method for multi-objective recommender systems, *Information Sciences* 572 (2021) 558–573.
- [33] K. Kapoor, V. Kumar, L. Terveen, J. A. Konstan, P. Schrater, " i like to explore sometimes" adapting to dynamic user novelty preferences, in: Proceedings of the 9th ACM Conference on Recommender Systems, 2015, pp. 19–26.
- [34] P. Li, M. Que, Z. Jiang, Y. Hu, A. Tuzhilin, Purs: Personalized unexpected recommender system for improving user satisfaction, *Fourteenth ACM Conference on Recommender Systems* (2020).
- [35] R. Mehrotra, C. Shah, B. Carterette, Investigating listeners’ responses to divergent recommendations, in: Proceedings of the 14th ACM Conference on Recommender Systems, RecSys ’20, Association for Computing Machinery, New York, NY, USA, 2020, p. 692–696. URL: <https://doi.org/10.1145/3383313.3418482>. doi:10.1145/3383313.3418482.
- [36] T. Qian, Y. Liang, Q. Li, X. Ma, K. Sun, Z. Peng, Intent disentanglement and feature self-supervision for novel recommendation, *ArXiv abs/2106.14388* (2021).

- [37] P. G. Campos, F. Díez, I. Cantador, Time-aware recommender systems: a comprehensive survey and analysis of existing evaluation protocols, *User Modeling and User-Adapted Interaction* 24 (2014) 67–119.
- [38] B. Veloso, B. Malheiro, J. C. Burguillo, J. Foss, Personalised fading for stream data, in: *Proceedings of the Symposium on Applied Computing*, 2017, pp. 870–872.
- [39] Y. Wu, K. Li, G. Zhao, X. Qian, Long-and short-term preference learning for next poi recommendation, in: *Proceedings of the 28th ACM international conference on information and knowledge management*, 2019, pp. 2301–2304.
- [40] Y. Gu, Z. Ding, S. Wang, D. Yin, Hierarchical user profiling for e-commerce recommender systems, in: *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 223–231.
- [41] J. Wang, R. Louca, D. Hu, C. Cellier, J. Caverlee, L. Hong, Time to shop for valentine’s day: Shopping occasions and sequential recommendation in e-commerce, in: *Proceedings of the 13th International Conference on Web Search and Data Mining*, 2020, pp. 645–653.
- [42] F. A. Chowdhury, Y. Liu, K. Saha, N. Vincent, L. Neves, N. Shah, M. W. Bos, Ceam: The effectiveness of cyclic and ephemeral attention models of user behavior on social platforms, in: *ICWSM*, 2021.
- [43] R. Verachtert, L. Michiels, B. Goethals, Are we forgetting something? correctly evaluate a recommender system with an optimal training window, in: *Proceedings of the Perspectives on the Evaluation of Recommender Systems Workshop*, volume 3228, 2022.
- [44] N. Lathia, S. Hailes, L. Capra, X. Amatriain, Temporal diversity in recommender systems, in: *Proceedings of the 33rd international ACM SIGIR conference on Research and development in information retrieval*, 2010, pp. 210–217.
- [45] M. Zaheer, A. Ahmed, A. J. Smola, Latent lstm allocation: Joint clustering and non-linear dynamic modeling of sequence data, in: *International Conference on Machine Learning*, PMLR, 2017, pp. 3967–3976.
- [46] M. M. Tanjim, C. Su, E. Benjamin, D. Hu, L. Hong, J. McAuley, Attentive sequential models of latent intent for next item recommendation, in: *Proceedings of The Web Conference 2020*, 2020, pp. 2528–2534.
- [47] C. Wang, W. Ma, M. Zhang, C. Chen, Y. Liu, S. Ma, Toward dynamic user intention: Temporal evolutionary effects of item relations in sequential recommendation, *ACM Transactions on Information Systems (TOIS)* 39 (2020) 1–33.
- [48] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, Ł. Kaiser, I. Polosukhin, Attention is all you need, *Advances in neural information processing systems* 30 (2017).
- [49] E. Fischer, D. Zoller, A. Hotho, Comparison of transformer-based sequential product recommendation models for the coveo data challenge (2021).
- [50] T. Lin, Y. Wang, X. Liu, X. Qiu, A survey of transformers, *arXiv preprint arXiv:2106.04554* (2021).
- [51] L. Li, L. Zheng, F. Yang, T. Li, Modeling and broadening temporal user interest in personalized news recommendation, *Expert Systems with Applications* 41 (2014) 3168–3177.
- [52] D.-T. Le, Y. Fang, H. W. Lauw, Modeling sequential preferences with dynamic user and context factors, in: *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, Springer, 2016, pp. 145–161.

A. Proofs

Proof of Lemma 3.1 If $\sigma^2(t) \geq \sigma^2(t-1) \geq 0$,

$$\begin{aligned} \frac{\sigma^2(t) - \sigma^2(t-1)}{\frac{1}{2}(\sigma^2(t) + \sigma^2(t-1))} &\leq \frac{\sigma^2(t)}{\frac{1}{2}(\sigma^2(t) + \sigma^2(t-1))} \\ &\leq \frac{\sigma^2(t)}{\frac{1}{2}\sigma^2(t)} = 2 \end{aligned}$$

If instead $\sigma^2(t) - \sigma^2(t-1) \leq 0$,

$$\begin{aligned} \frac{\sigma^2(t) - \sigma^2(t-1)}{\frac{1}{2}(\sigma^2(t) + \sigma^2(t-1))} &\geq -\frac{\sigma^2(t-1)}{\frac{1}{2}(\sigma^2(t) + \sigma^2(t-1))} \\ &\geq -\frac{\sigma^2(t-1)}{\frac{1}{2}\sigma^2(t-1)} = -2 \end{aligned}$$

The two extreme values are in fact achieved by $\sigma^2(t) = 0$ and $\sigma^2(t-1) = 0$ respectively; as we discard any such cases as outliers, the range is strictly speaking $(-2, 2)$. \square

Proof of Lemma 3.2

If $\sigma_j^2 = \bar{\sigma}_j^2$, $\overline{DispAdd}_j = DispAdd_j$ by construction.

In the reverse direction, if $|\overline{DispAdd}(t) - DispAdd(t)| = 0$,

$$\begin{aligned} 0 &= \frac{\bar{\sigma}^2(t) - \sigma^2(t-1)}{\frac{1}{2}(\bar{\sigma}^2(t) + \sigma^2(t-1))} - \frac{\sigma^2(t) - \sigma^2(t-1)}{\frac{1}{2}(\sigma^2(t) + \sigma^2(t-1))} \\ &= \frac{2(\bar{\sigma}^2(t) - \sigma^2(t))\sigma^2(t-1)}{(\bar{\sigma}^2(t) + \sigma^2(t-1))(\sigma^2(t) + \sigma^2(t-1))} \end{aligned}$$

Note that *Dispersion* is nonnegative by construction. Then if the denominator is nonzero, the numerator is zero only if $\bar{\sigma}^2(t) - \sigma^2(t) = 0$ given $\sigma^2(t-1) \neq 0$. \square