12-10-2023

# Warming Up the Cold Start: A Multimodal Approach

Anat Goldstein
*Ariel University*, anatgo@ariel.ac.il

Chen Hajaj
*Ariel University*, Chenha@ariel.ac.il

# Warming Up the Cold Start: A Multimodal Approach

*Research-in-Progress*

**Anat Goldstein**

Ariel University

anatgo@ariel.ac.il

**Chen Hajaj**

Ariel University

Chenha@ariel.ac.il

## Abstract

Recommendation systems are an important part of our daily lives, but they can struggle in cold-start scenarios, when there are no historical transaction data for items and users. This study addresses such cold-start scenarios. Our proposed two-phase solution involves first, identifying for each new item its most similar existing items using various item embeddings (image, textual-description, and attributes-based), and second, recommending frequently co-purchased items, alongside an augmented nearest neighbor collaborative filtering method. The proposed solution was evaluated using data from H&M's online store. Our preliminary outcomes highlight the effectiveness of textual-description embeddings for generating valuable recommendations under cold-start conditions. Notably, using textual embeddings, we accurately recommend at least one item among the top co-purchased items for 20% of new items, achieving substantial performance improvements over baseline strategies. This study contributes a versatile solution requiring no retraining of models and catering to scenarios with limited historical activity on websites.

**Keywords**

Recommendation systems, cold-start, embeddings, multimodal, collaborative filtering, product-network.

## Introduction

Recommendation systems have become an integral part of our daily lives. From online shopping to streaming services, we rely on recommendation systems to make personalized suggestions that match our preferences.
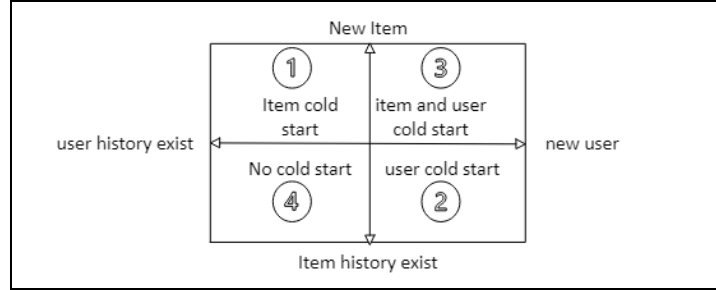
There are two primary approaches for generating personalized recommendations: The first and most common approach is collaborative filtering (CF), which recommends based on past interactions that various users have had with different items (e.g., purchases they've made or ratings they've given). According to this approach, each user is characterized by their interactions with the items presented on the website, and each item is characterized by the interactions that different users have had with it (e.g., ratings provided by various users or the presence of previous purchases).

Two strategies can be employed when recommending items to a user: First, identifying users with similar past interactions to the target user and suggesting items they have purchased (or given high ratings to) that the user has not yet purchased. Second, identifying items similar to those the user had engaged with in the past and recommending them. It's worth noting that there are different implementations of CF, including techniques like nearest neighbors (Nikolakopoulos et al. 2021), matrix factorization (Koren et al. 2009; Li et al. 2008), and deep learning techniques (Barrett et al. 2017).

The second recommendation approach is content-based filtering (CB). This approach (Lops et al. 2010) leverages the attributes of items (e.g., visual, functional, and technical attributes) and aligns them with user profiles and demographic information.

While the first approach is more popular and generally considered more effective, it does have a drawback: when an item or user is new to the system, their history of interactions is unknown, making it difficult to generate effective recommendations for them. This is known as the cold start problem. Figure 1 illustrates possible cold start scenarios. Specifically, it shows four scenarios: (1) *item cold start* – when we recommend a new item to existing users (i.e., users with history of interactions), (2) *user cold start* – when we

recommend existing items (e.g., items with past interactions) to a new user, (3) *item and user cold start* – when we recommend new items to new users, and (4) *no cold start* - when we recommend existing items to existing users, that is, both items and users have history of interactions.



**Figure 1. Cases in which different recommendation methods should be used**

To address the cold start problem, different hybrid methods were proposed in the literature, that combine CB recommendations and CF, when interaction history is missing. Combining these two approaches allows the system to make more accurate recommendations, even for new users and items. Nevertheless, despite the extensive literature that focuses on the cold-start problem (Cai et al. 2023; Heidari et al. 2022; Jafri et al. 2022; Lam et al. 2008; Patro et al. 2023; Zhang et al. 2013), it remains an ongoing research challenge, particularly in scenarios of the third type - item and user cold start. Our research is specifically focused on addressing this type of scenario.

We propose a two-phase CB solution, specifically designed for scenarios lacking any user-item interaction history. The first phase involves generating new item embeddings based on different item-modalities: images, textual descriptions, and item attributes. Furthermore, we compare these diverse embedding types and assess their effectiveness. Subsequently, we identify the most similar existing item to the new item by comparing their embeddings.

In the second phase, utilizing a predefined *item network*, we identify, from the existing items, those that are frequently purchased in conjunction with the item found in the first phase. As elucidated below, while the second phase doesn't necessitate historical user interaction data, we also propose a variation of the second phase that incorporates users' interaction history if available.

To test the proposed method, we use data from H&M's online fashion shop, which is publicly available on Kaggle.com[1]. The realm of fashion presents numerous challenges for recommendation systems (Kula, 2015): Firstly, fashion stores typically offer a vast array of items, resulting in sparse interaction data. Secondly, in this particular context, the most pertinent, fashionable, items are the newly released products, contributing to the *item cold start* scenario. Thirdly, a substantial number of users remain unidentifiable or are first-time users, giving rise to the user cold start scenario.

Using the H&M dataset, we remove a random subset of items (10%) and their respective transactions to simulate cold start scenarios. We use our recommendation method to recommend items to users in the test set. We then evaluate the recommendations against the actual transactions using common metrics.

Our preliminary results are promising. We show that the textual embeddings of the new items are more informative and effective for recommendation compared to image embeddings, compared to attribute embeddings and even compared to a multimodal representation of image and textual description. Employing our approach with textual embeddings in an *item and user cold start* scenario, we are able to accurately pinpoint at least one item out of the actual five most co-purchased items for 40% of the new items. Our predicted co-purchased items obtain a recall of 12.6%, precision of 13.4 % and Jaccard index of 8% (k=5 items). This marks a substantial advancement compared to the baseline recommendation strategy of the five most popular items, where the recall stands at a mere 0.4%, the precision at 0.5%, and the Jaccard index at 0.1%.

---

[1] https://www.kaggle.com/competitions/h-and-m-personalized-fashion-recommendations/

# Data

We utilize a publicly available dataset extracted from H&M's online store. This dataset encompasses the purchasing history of customers spanning over time, along with comprehensive item details, including images, textual descriptions, and supplementary attributes such as item type, department, and color. In total, the dataset includes 104,500 different products, 103,426 of which include image and textual description. Additionally, our dataset comprises pertinent customer information on 1,371,980 customers, including customer ID, age, postal code, frequency of fashion news consumption, club membership status, and customer activity. Finally, we have data on 31, 788,424 transactions. This transactional data includes essential details such as customer ID, article ID, price, and the transaction channel.

# Research Method
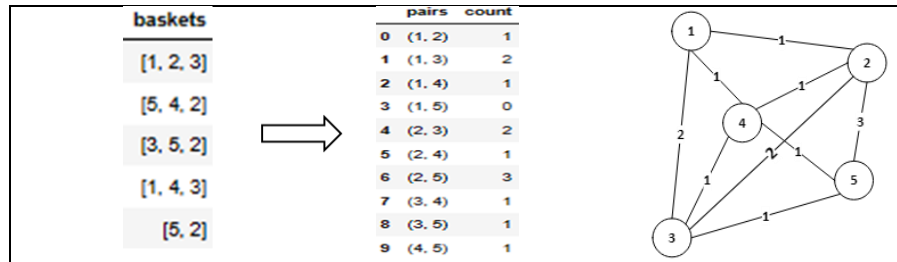
## *Creating multi-modal embeddings*

For each item, we create different embeddings that capture different item modalities: image-based, text-based, and attribute-based. These embeddings are compact numerical representations that capture the underlying semantic relationships between items.

To create image-based item embeddings, we employed the pretrained ResNet model (He et al. 2016), which transforms images into meaningful numerical vectors. To create textual description-based embedding we used the pretrained RoBERTa language model (Liu et al. 2019), which translates items' textual descriptions into dense numerical vectors that preserve semantic meaning. Finally, to create attribute-based embeddings, we used a vector of values that encapsulates various predefined features of the item. Most of these features were categorical in nature (e.g., product type, color, graphical appearance, and department) and were represented using one-hot encodings.

## *Creating an item network based on previous transactions*

Using the transaction history from the website, we construct an item network, also referred to as a product network. This network is bidirectional in nature, with nodes representing items and edges connecting nodes indicating instances of co-purchases. The weight assigned to each edge corresponds to the frequency with which the connected items were bought together by the same user during a single visit to the website. Such an item network is commonly employed in academic literature to capture latent relationships between products (Goldstein et al. 2022; Oestreicher-Singer and Sundararajan 2012).

Figure 2 illustrates how the item network is generated based on user baskets containing co-purchased items. Assuming that we have five items (labeled as 1 to 5) that were purchased in five distinct user baskets, as depicted on the left side of Figure 2, the resultant network is presented on the right side of the figure. This presentation is provided in both tabular form (indicating the frequency of joint purchases for each pair of items) and as a graphical representation.



**Figure 2. Creating item network from user baskets**

## *Recommending new items based on similarity*

We present a two-phase solution, depicted in Figure 3. In the first phase, given a new item, denoted as $item_i$, we generate an embedding for the item as discussed earlier (step 1a in Figure 3). Subsequently, we determine the most similar item to $item_i$ by computing the cosine similarity between the embeddings of the items (step 1b in Figure 3).

Assuming that the most similar product is $item_j$, we move to the second phase. In this phase, we utilize the pre-existing item network to identify items frequently purchased in conjunction with $item_j$, thereby making recommendations (step 2a in Figure 3).

We explore four variations of this procedure, each encompassing distinct embeddings (image, textual description, attribute-based, and multimodal representation which combines both image and description) and four diverse recommendation methods:

*Method 1*: Identifying the item most similar to the new item *item_i* and determining the k most commonly purchased items alongside it (depicted in Figure 4A).

*Method 2*: Identifying the k items most similar to a new item *item_i* and determining the most frequently purchased item in conjunction with each of the k items (depicted in Figure 4B).

*Method 3*: Identifying the k items most similar to a new item *item_i* and recommending them.

Method 4: combining methods 2 and 3: we recommend m items based on method 3 (m most similar items to $item_i$) and k-m items based on method 2 (items most purchased with the most similar items to $item_i$).

Additionally, we offer nearest-neighbor collaborative filtering (step 2b in Figure 3) recommendations by augmenting an 'artificial' vector of interactions for the new item. This vector replicates the interaction vector of $item_j$. This inclusion allows the recommendation system to suggest the new item to existing users as well.
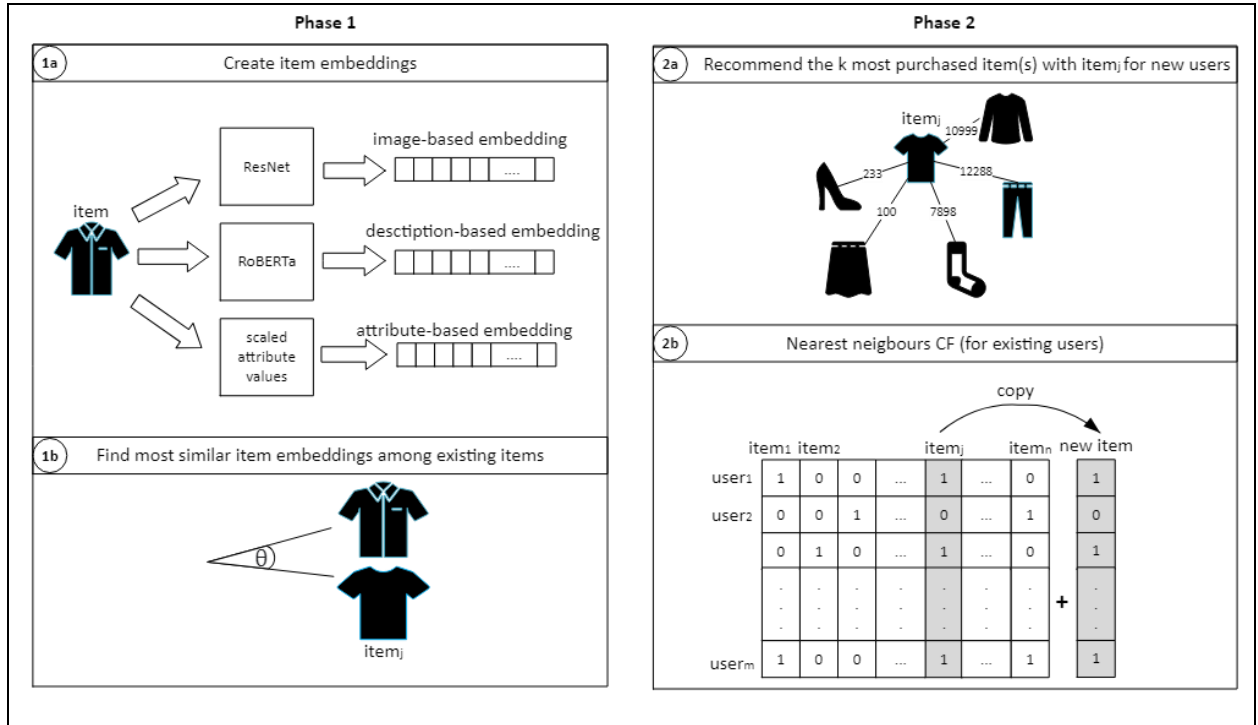


**Figure 3. Proposed recommendation method steps**
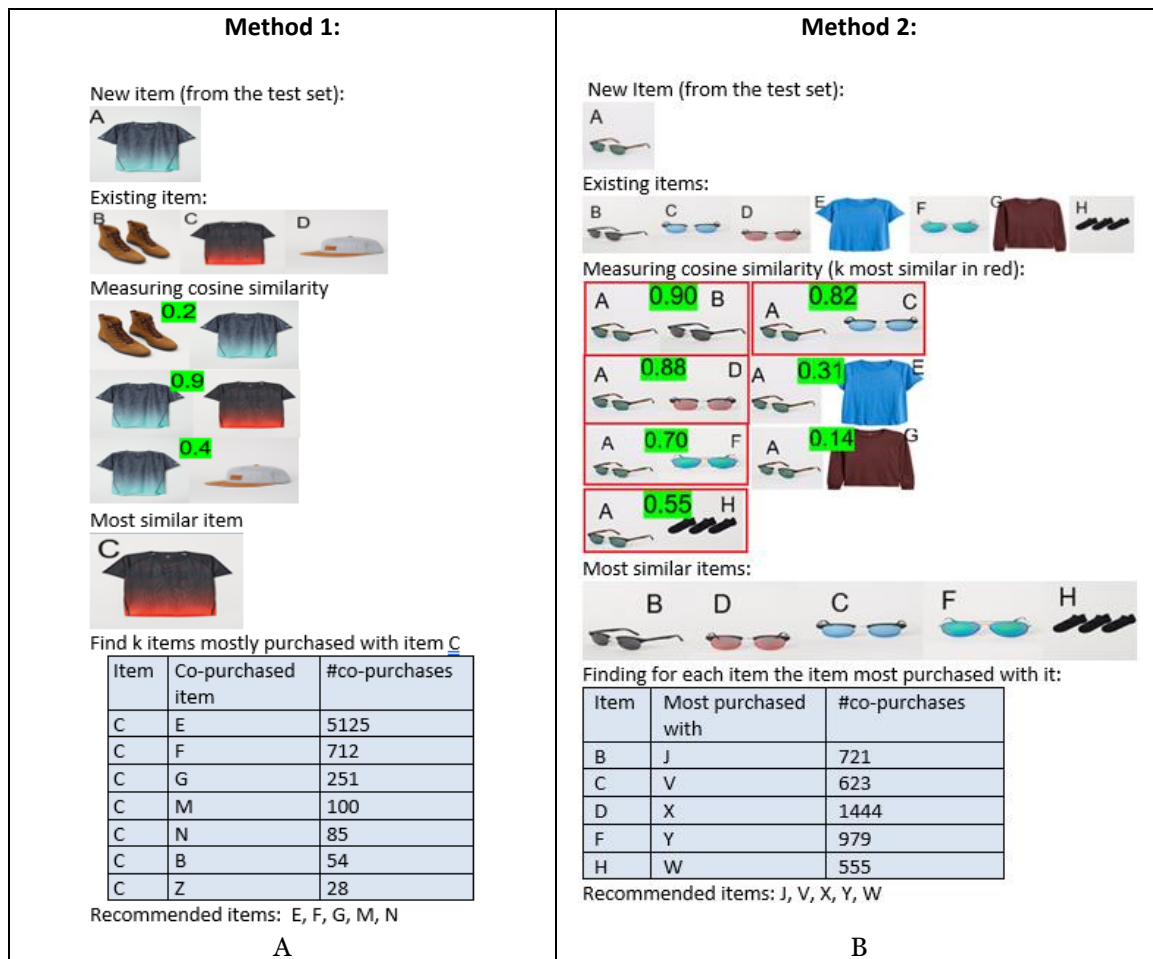
## *Evaluating model recommendation*

We implement two evaluation procedures. The first evaluation centers around the initial phase of the proposed solution and assesses the effectiveness of the different item embeddings used for finding the most

similar existing items to the new items. The second evaluation is geared towards assessing the caliber of user recommendations.

**Evaluation of Embeddings Effectiveness (First Phase)**

We randomly partition the items within the dataset into two categories, namely 'existing items' (constituting 90% of the data) and 'new items' (representing 10% of the data). The 'new items' are reserved for testing the model, while the transactions involving 'existing items' are used for establishing the item network, which was previously discussed.

Additionally, we construct a complete item-network from the entire transaction data, encompassing all items. This network serves as the 'ground truth' against which we evaluate the recommendations, as explained below.



**Figure 4. Examples of the first two recommendation methods: On the left side (A), we demonstrate the first recommendation method and, on the right (B), the second method. The third method is actually the first three steps on B.**

To gain preliminary insights into the efficacy of the proposed method, we begin by determining whether the recommended items corresponding to each new item genuinely represent the most frequently co-purchased items, according to the ground truth, which comprises the complete transaction history. Consequently, for each 'new' item, we extract the top five recommended items (as per their similarity to the existing items and their co-purchases within the existing item network). We then cross-reference these results with the complete item network, encompassing both new and existing items, and identify the five most co-purchased items for each new item. This enables us to compare the extent to which the

recommended items align with the five most co-purchased items, as per the ground truth. To conduct this comparison, we employ four metrics:

1. *Hit Rate* – Signifying the percentage of new items for which the model accurately predicted (recommended) at least one co-purchased item out of the five recommended items (k = 5).
2. *Jaccard Index* – Representing the intersection of recommended items and relevant items, divided by the union of recommended items and relevant items.
3. *Recall@k* – Indicating the ratio of relevant k recommended items out of the collection of top k relevant items (i.e., the actual top five co-purchased items).
4. *Precision@k* – Conveying the proportion of relevant recommended items out of the k recommended items.

We repeat this process 5 times by randomly dividing the data into 'new' vs. 'existing' items, providing both mean values and standard deviations. This preliminary evaluation allows us to validate the efficacy of the item embeddings utilized to identify the most similar items in the first phase of our recommendation solution.

We proceed to compare the results yielded by our approach against two baseline methods:

(1) A straightforward (naïve) model that recommends the five top items (most popular items).
(2) A more sophisticated variant of the naïve model, which recommends the five most popular items from the same product group as the item in question (e.g., bras, socks, dresses)

**User Recommendation Evaluation**

As mentioned above, we partition the items in the dataset into: 'existing items' (comprising 90% of the data) and 'new items' (constituting 10% of the data). All transactions or interactions involving the new items are extracted from the training set and incorporated into the test set. This creates a scenario in which the recommendation model must generate suggestions from a pool of items lacking any interaction data.

For every user present in the test set, we provide recommendations for five items (using both the item-network and our augmented nearest-neighbor collaborative filtering mentioned above) and subsequently assess the recall@k, precision@k, and the Jaccard index (with k set to 5) across all users.

To establish a basis of comparison, we contrast our outcomes with those achieved by a baseline nearest neighbor collaborative filtering algorithm.

This process is reiterated by randomly dividing the data into 'new' vs. 'existing' items five times, allowing us to present both mean values and standard deviations of the model performance metrics.

# Preliminary Results

Up to this point, we have completed the initial evaluation procedure, concentrating on the effectiveness of diverse embedding approaches, as quantified by the metrics: hit-rate, recall@k, precision@k, and the Jaccard index. For our initial investigation, we employed values of k equal to 5 and 10, as is commonly found in the literature (Adomavicius and Zhang 2016; Cai et al. 2023; Patro et al. 2023; Yang et al. 2012). we are presenting the results for a 5-item recommendation here. It's important to note that recommending 10 items yielded results of a similar nature and magnitude.

Table 1 presents our promising preliminary findings in an item and user cold start scenario. The results reveal that embeddings based on item textual descriptions displayed superior effectiveness (highlighted in bold) compared to those derived from images, predefined attributes, and even a combined (multimodal) representation of textual descriptions and images. The only exception to this trend is the hit-rate metric, which was slightly higher in the multimodal-based recommendation. Additionally, these description-based embeddings significantly outperformed baseline strategies that involve recommending popular items. Method 4, in which the recommendation set includes both similar items to the new item and their co-purchased items, consistently produced the best results for every type of item representation (textual description-based, image-based, attribute-based, and multimodal).

Notably, when utilizing our fourth recommendation method in conjunction with text-based embeddings in situations where neither item nor user interaction history is accessible, we achieve a commendable accuracy

in recommending at least one item (out of five) that aligns with the actual top five co-purchased items for 40% of the new items. Our evaluation of the first phase using Method 4 reveals a recall@5 of 12.6%, precision@5 of 13.4%  and a Jaccard index of 8%.

| Approach | Method | Hit rate (>=1) | Recall@5 | Precision@5 | Jaccard index |
|---|---|---|---|---|---|
| Baseline | Popularity | 0.005 (0.005) | 0.001 (0.001) | 0.001 (0.001) | 0.001 (0.001) |
| Advanced baseline | Popularity in product group | 0.022 (0.001) | 0.005 (0.0003) | 0.004(0.0003) | 0.003 (0.0001) |
| Attribute-based | Method 1 | 0.017 (0.007) | 0.004 (0.002) | 0.004 (0.002) | 0.002 (0.002) |
| | Method 2 | 0.031 (0.002) | 0.006 (0.005) | 0.007 (0.006) | 0.004 (0.003) |
| | Method 3 | 0.025 (0.002) | 0.005 (0.005) | 0.005 (0.005) | 0.003 (0.003) |
| | Method 4 | 0.031 (0.002) | 0.007 (0.004) | 0.006 (0.003) | 0.004 (0.002) |
| Image-based | Method 1 | 0.107 (0.003) | 0.033 (0.008) | 0.033 (0.008) | 0.021 (0.005) |
| | Method 2 | 0.118 (0.002) | 0.027 (0.006) | 0.029 (0.001) | 0.016 (0.004) |
| | Method 3 | 0.144 (0.005) | 0.032 (0.001) | 0.032 (0.001) | 0.019 (0.007) |
| | Method 4 | 0.125 (0.004) | 0.033 (0.001) | 0.029 (0.001) | 0.018 (0.006) |
| Description-based | Method 1 | 0.261 (0.004) | 0.082 (0.001) | 0.082 (0.001) | 0.052 (0.009) |
| | Method 2 | 0.271 (0.004) | 0.078 (0.001) | 0.097 (0.002) | 0.053 (0.008) |
| | Method 3 | 0.388 (0.005) | 0.12 (0.002) | 0.12 (0.002) | 0.076 (0.001) |
| | Method 4 | 0.398 (0.005) | **0.126 (0.002)** | **0.134 (0.002)** | **0.082 (0.001)** |
| Multimodal | Method 1 | 0.252 (0.004) | 0.084 (0.001) | 0.084 (0.009) | 0.054 (0.006) |
| | Method 2 | 0.267 (0.004) | 0.073 (0.001) | 0.089 (0.001) | 0.049 (0.007) |
| | Method 3 | 0.383 (0.004) | 0.108 (0.001) | 0.108 (0.001) | 0.066 (0.009) |
| | Method 4 | **0.404 (0.005)** | 0.124 (0.002) | 0.118 (0.002) | 0.074 (0.001) |

**Table 1. Item embedding performance evaluation.**

## Conclusion

Despite the extensive literature addressing the cold start problem, it remains a challenge for online stores in general, and particularly for online fashion stores, where the most relevant items are the newest releases, which lack purchase history and where many users are occasional customers who lack transaction history. In the study presented in this article, we propose methods tailored for scenarios where no historical data is available for items or users. These methods rely on measuring similarity between new and existing items and identifying products that are likely to be purchased alongside the existing ones using an item-network. Our initial results demonstrate the effectiveness of various similarity measurement methods based on items' textual description embeddings, image embeddings and attribute embeddings. These results indicate that measuring similarity using the textual description of the item achieves the best performance. Additionally, our findings suggest that recommending a combination of similar items and their co-purchased items is the most effective method.

It's worth noting that recommendations based on text-based embeddings outperformed recommendations based on multimodal item representation, which combines textual descriptions and images of the item. This suggests that adding image information to the item representation did not improve the recommendations. One potential reason for this could be that we used a generic pretrained model rather than an image model fine-tuned to the context of fashion items. In subsequent phases of our research, we will explore whether using fine-tuned image embeddings can enhance recommendations.

Furthermore, in the subsequent phases of our research, we will evaluate the recommendations provided by the proposed method across various configurations and cold-start scenarios, comparing them against established recommendation methods. We will also explore additional methods that employ graph structure embeddings.

The contribution of our method is twofold: firstly, it does not require retraining the recommendation model to include new items or users, as do many of the existing hybrid approaches; Secondly, the method can be applied in any context, and is particularly suited for scenarios where there is no historical activity on the website for both products and users and in contexts where collaborative filtering has been shown to be less suitable, such as in the tourism context (Jannach et al. 2007).

# References

Adomavicius, G., and Zhang, J. 2016. "Classification, Ranking, and Top-K Stability of Recommendation Algorithms," *INFORMS Journal on Computing* (28:1), pp. 129–147. (https://doi.org/10.1287/ijoc.2015.0662).

Barrett, R., Cummings, R., Agichtein, E., Gabrilovich, E., He, X., Liao, L., Zhang, H., Nie, L., Hu, X., and Chua, T.-S. 2017. "Neural Collaborative Filtering," *Proceedings of the 26th International Conference on World Wide Web*, pp. 173–182. (https://doi.org/10.1145/3038912.3052569).

Cai, D., Qian, S., Fang, Q., Hu, J., and Xu, C. 2023. "User Cold-Start Recommendation via Inductive Heterogeneous Graph Neural Network," *ACM Transactions on Information Systems* (41:3), pp. 1–27. (https://doi.org/10.1145/3560487).

Goldstein, A., Oestreicher-Singer, G., Barzilay, O., and Yahav, I. 2022. "Are We There Yet? Analyzing Progress in the Conversion Funnel Using the Diversity of Searched Products," *MIS Quarterly* (46:4), pp. 2015–2054.

He, K., Zhang, X., Ren, S., and Sun, J. 2016. "Deep Residual Learning for Image Recognition," in *Computer Vision and Pattern Recognition (CVPR), 2016*.

Heidari, N., Moradi, P., and Koochari, A. 2022. "An Attention-Based Deep Learning Method for Solving the Cold-Start and Sparsity Issues of Recommender Systems," *Knowledge-Based Systems* (256), p. 109835. (https://doi.org/10.1016/j.knosys.2022.109835).

Jafri, S. I. H., Ghazali, R., Javid, I., Mahmood, Z., and Hassan, A. A. A. 2022. "Deep Transfer Learning with Multimodal Embedding to Tackle Cold-Start and Sparsity Issues in Recommendation System," *PLoS ONE* (17:8), p. e0273486. (https://doi.org/10.1371/journal.pone.0273486).

Jannach, D., A, M. Z., A, M. J., Seidler, O., and Warmbad-villach, C. T. 2007. "Developing a Conversational Travel Advisor with ADVISOR SUITE," in *Information and Communication Technologies in Tourism 2007*, Sigala, Mich, and Murphy (eds.), pp. 43–52. (https://doi.org/10.1007/978-3-211-69566-1_5).

Koren, Y., Bell, R., and Volinsky, C. 2009. "Matrix Factorization Techniques for Recommender Systems," *Computer* (42:8), pp. 30–37. (https://doi.org/10.1109/mc.2009.263).

Lam, X. N., Vu, T., Le, T. D., and Duong, A. D. 2008. "Addressing Cold-Start Problem in Recommendation Systems," *Proceedings of the 2nd International Conference on Ubiquitous Information Management and Communication*, pp. 208–211. (https://doi.org/10.1145/1352793.1352837).

Li, Y., Liu, B., Sarawagi, S., and Koren, Y. 2008. "Factorization Meets the Neighborhood: A Multifaceted Collaborative Filtering Model," *Proceedings of the 14th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, pp. 426–434. (https://doi.org/10.1145/1401890.1401944).

Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer, L., and Stoyanov, V. 2019. "RoBERTa: A Robustly Optimized BERT Pretraining Approach," *ArXiv*.

Lops, P., Gemmis, M. de, and Semeraro, G. 2010. *Recommender Systems Handbook*, pp. 73–105. (https://doi.org/10.1007/978-0-387-85820-3_3).

Nikolakopoulos, A. N., Ning, X., Desrosiers, C., and Karypis, G. 2021. *Recommender Systems Handbook*, pp. 39–89. (https://doi.org/10.1007/978-1-0716-2197-4_2).

Oestreicher-Singer, G., and Sundararajan, A. 2012. "Recommendation Networks and the Long Tail of Electronic Commerce," *MIS Quarterly* (36:1), pp. 65–83.

Patro, S. G. K., Mishra, B. K., Panda, S. K., Kumar, R., Long, H. V., and Taniar, D. 2023. "Cold Start Aware Hybrid Recommender System Approach for E-Commerce Users," *Soft Computing* (27:4), pp. 2071–2091. (https://doi.org/10.1007/s00500-022-07378-0).

Yang, X., Steck, H., Guo, Y., and Liu, Y. 2012. "On Top-k Recommendation Using Social Networks," *Proceedings of the Sixth ACM Conference on Recommender Systems*, pp. 67–74. (https://doi.org/10.1145/2365952.2365969).

Zhang, D., Hsu, C.-H., Chen, M., Chen, Q., Xiong, N., and Lloret, J. 2013. "Cold-Start Recommendation Using Bi-Clustering and Fusion for Large-Scale Social Recommender Systems," *IEEE Transactions on Emerging Topics in Computing* (2:2), pp. 239–250. (https://doi.org/10.1109/tetc.2013.2283233).