

E-commerce products

Project Report for NLP Course, Winter 2023

S. Rećko
WUT

01151399@pw.edu.pl

M. Sperkowski
WUT

01151430@pw.edu.pl

P. Tomaszewski
WUT

01151442@pw.edu.pl

K. Ułasiak
WUT

01151444@pw.edu.pl

supervisor: A. Wróblewska
WUT

anna.wroblewska1@pw.edu.pl

Abstract

The rapid growth of e-commerce, fueled by shifts in consumer behavior and accentuated by the COVID-19 pandemic, has underscored the need for efficient recommendation systems. This article presents a novel framework that leverages machine learning and NLP techniques to develop an automated and comprehensive recommendation system for e-commerce platforms. The focus lies on accurately measuring product similarity based on various attributes, including taxonomy, textual descriptions, and titles. The framework, designed for real-time calculations, aims to enhance e-commerce search functionalities and improve the overall user experience. The project employs state-of-the-art language models such as BERT, DistilBERT, and RoBERTa, incorporating a custom metric to evaluate the performance of these models. The experiments include training procedures with a Triple Loss function, and the results demonstrate promising advancements in e-commerce recommendation systems.

1 Introduction

The rise of e-commerce has been an unprecedented change in the sales market. Over the past years, e-commerce has witnessed exponential growth, accelerated further by global shifts in consumer behavior, especially the increased adoption of online shopping. The convenience, accessibility, and wide array of products available online have transformed how people shop. The COVID-19 pandemic further expedited this shift, with many consumers transitioning to online platforms for everyday needs. As e-commerce contin-

ues to expand, the need for efficient recommendation systems becomes even more critical, (Zhang et al., 2019). With an ever-growing number of products available, these systems play a pivotal role in helping consumers navigate through the vast array of choices, making their shopping experiences more personalized, streamlined, and enjoyable (Bobadilla et al., 2013). They do so by harnessing the power of data, algorithms, and user behavior to offer personalized product recommendations.

The overarching goal of this project is development of an innovative and robust framework leveraging machine learning and NLP techniques, to create an automated and comprehensive recommendation system. The architecture is capable of measuring the similarity and diving taking into account the different levels of dissimilarity between products within a e-commerce platforms. This research aims to achieve an advanced level of accuracy in defining and establishing similarity metrics between products based on various attributes, including taxonomy (categories), textual descriptions and titles. To be practical for the end users, calculations must be performed quickly. Therefore, the calculations should work in real-time, so a limit of at least 0.01 seconds (10ms) per pair of products is set. This is considered an important requirement for the project.

The primary scientific objective is to construct a sophisticated similarity measurement that accurately captures the semantics and nuances of product relationships, distinguishing between identical, closely related (such as variations within the same product line differing in specifications like memory sizes), and entirely distinct products. This entails the exploration and development of algorithms that can comprehend and differentiate the diverse relationships among products, accommodating subtle variations in features while still recognizing commonalities.

Furthermore, the project will focus on innovating techniques for automatically extracting key information from item descriptions and titles using Large Language Models. These generated attributes will supplement existing data, enriching the product information available on e-commerce platforms and possibly creating a better representation of the products.

Ultimately, the scientific aim is to significantly enhance e-commerce search functionalities, facilitate the simple yet efficient introduction of new products into online marketplaces, and provide users with access to a comprehensive range of available products and complementary items. This research endeavour will contribute to advancing the capabilities of e-commerce platforms by refining the accuracy and depth of product recommendations, thereby improving the overall user experience in online product search and comparison.

1.1 Research questions

In e-commerce, the sheer volume and diversity of products pose a significant challenge in accurately measuring their similarity across multiple dimensions. The current methodologies for comparing and categorizing products often fall short in capturing the nuanced differences and similarities between items. This leads to sub-optimal search experiences for users and hampers the efficiency of introducing new products into e-commerce platforms. The need for an automated and precise system to measure product similarity, considering varying attributes and features, is crucial for enhancing product search accuracy and user satisfaction.

Research Questions and Hypothesis:

- Q: Is there a measure that can incorporate the taxonomy of products in a similarity measurements and if so, can it be incorporated to the loss function for fine-tuning transformer models?
- H: Such measure can be found, however the incorporation may prove difficult. It could require much more inputs to the model which would significantly increase the processing time.
- Q: Can we leverage the Large Language Models complexity to augment the data and achieve better representation of the products?

H: Generative power and generalisation of LLMs can be used to augment the dataset, as they prove to be able to solve many real world problems. Issue might arise concerning the API limits for free LLMs or limited size of these models.

Q: Can real-time performance of the pair similarity measurement be achieved using transformer models? (Reaching 10ms per pair comparison)

H: Such performance is achievable, however it might require novel approaches to minimize the inference time. The complications might come from limited computational power available in the project.

Q: Can the approach from one dataset be extended to another, possibly in other language and can it achieve similar results?

H: The architecture should be carefully developed to solve a similar task with minimal changes. However, the performance of the models is heavily dependent on the amount of available data, which is typically harder to come by in languages other than english.

1.2 Report structure

In Section 2 we describe the scientific literature related to our project. We focus on the state-of-the-art in NLP and text embeddings, e-commerce recommendation systems, similarity measures and Large Language Models. Furthermore, the closely related works are described in detail, together with the datasets used in research. In Section 3 the research methodology adopted in this project, especially the techniques and tools for both research and result analysis. In Section 4 the experiments on the WDC datasets, including the exploratory data analysis, data augmentation using LLMs and finetuning of the transformer model. In Section 5 a discussion of our results in comparison to the literature is made. In Section 6 the project conclusions are written, summing up what has been done and proposing future works for this project and the domain.

2 Related works

Machine learning approaches for extracting meaningful attributes from product descriptions include NLP techniques, word embeddings, and Named

Entity Recognition (NER). These attributes can be integrated into the similarity measurement process by converting them into product feature vectors. Utilizing metrics such as cosine similarity or Euclidean distance on these vectors allows for an effective quantification of product similarity, enhancing the overall recommendation system.

Siamese Neural Networks (Koch et al., 2015), Graph Neural Networks (Scarselli et al., 2008) (GNN), and Transformer models (Vaswani et al., 2017) like BERT (Devlin et al., 2019) excel in capturing semantics and nuanced relationships between products. Siamese networks are adept at understanding subtle differences, GNNs model complex dependencies and Transformers provide contextualized embeddings, collectively offering a robust framework for differentiating between identical, slightly different, and entirely distinct items.

Integrating generated product attributes into existing e-commerce platforms can be achieved by developing an attribute-based search functionality, enhancing recommendation systems, and optimizing the introduction of new products. By leveraging these attributes, platforms can offer more personalized search options, improve recommendation accuracy, and streamline the process of introducing new products efficiently into the market, ultimately enhancing the overall user experience.

Various strategies can be employed to reach minimal pair comparison time while utilizing state-of-the-art (SOTA) models with complex architectures. Techniques like model quantization (Polino et al., 2018) and pruning (Liu et al., 2018) help reduce the computational load, while hardware acceleration using specialized processors like GPUs or TPUs speeds up inference. Additionally, caching and batch processing can be implemented to precompute certain calculations and perform parallelized comparisons, ensuring efficient and real-time processing without compromising the sophistication of the underlying models. Distilling the knowledge of a bigger model to a smaller one, by enforcing similar outputs on the training data is another method commonly used for creating smaller networks (Gou et al., 2021).

In recent years, Large Language Models (LLMs) have emerged as transformative tools across various domains, showcasing their remarkable versatility and utility. These models, such as GPT-3 (Brown et al., 2020), have demonstrated an unprecedented ability to understand and generate

human-like text, enabling advancements in natural language processing, text generation, and information retrieval. In artificial intelligence, LLMs have been pivotal in enhancing chatbot capabilities, language translation, and content creation. Moreover, they have proven instrumental in automating mundane tasks, facilitating more efficient data analysis, and even contributing to the development of novel applications in healthcare (Thirunavukarasu et al., 2023), finance (Wu et al., 2023), and education (Kasneci et al., 2023). The extensive capabilities of Large Language Models (Wei et al., 2022) underscore their potential to revolutionize how we interact with technology, opening up new frontiers for innovation and problem-solving across diverse fields.

2.1 Available Datasets

In recent years, entity resolution has shifted towards deep learning-based matching methods, necessitating large training data. Traditional benchmark datasets often prove inadequate for evaluating these methods due to their limited size and source diversity. The "Web Data Commons" (WDC) dataset (Web Data Commons - Training Dataset and Gold Standard for Large-Scale Product Matching) addresses these challenges by offering a substantial volume of data, including 16 million English-language offers, sourced from 79 thousand websites. This diversity and scale make it an ideal choice for assessing deep learning-based matches and improving their evaluation and comparison. The dataset includes categorization using distant supervision from Amazon product data. Lexica containing terms and their TF-IDF scores for 26 product categories (look Figure 1) were created using publicly available Amazon product reviews and metadata. Each offer in the dataset is assigned the product category whose terms maximize the sum of overlapping TF-IDF scores. In cases with minimal overlap, the offer is categorized as "not found". We exclusively utilized the "Gold Standard" for the training of our product matching method. The gold standard, derived from the English product data corpus, comprises a set of 1,100 pairs of offers from each of the four product categories: Computers Accessories, Camera & Photo, Watches, and Shoes. For each product, the gold standard includes two matching pairs of offers (positives) and five or six non-matching pairs of offers (negatives).

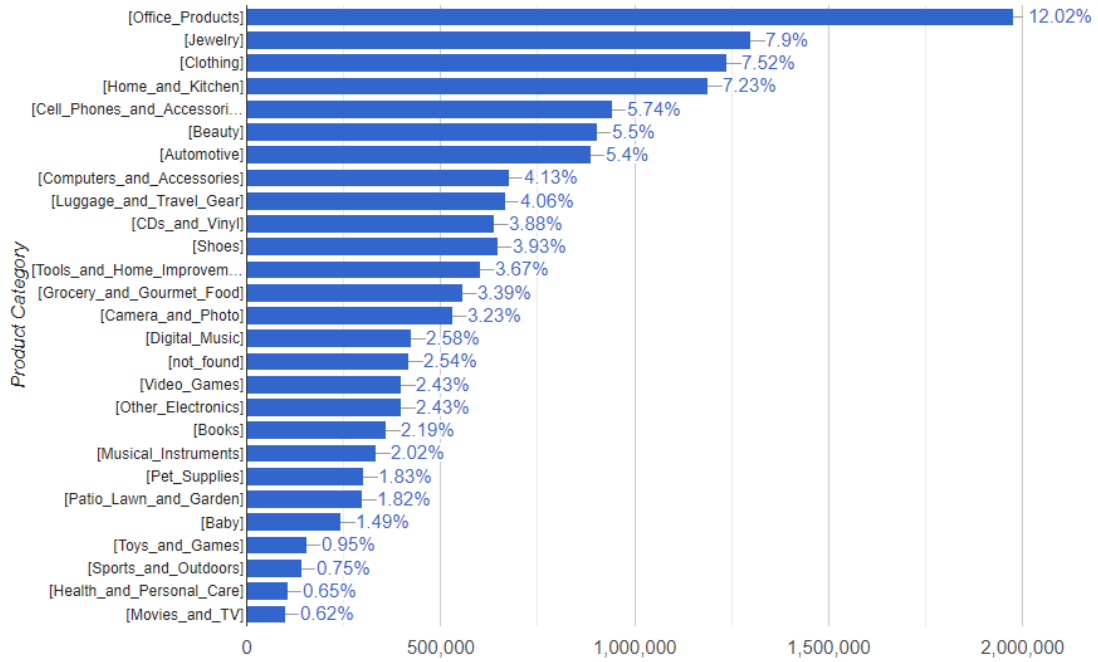


Figure 1: Distribution of offer entities per category in the English Training Set.

The Polish dataset for product matching created was created by the authors of (Michał Mozdzzonek, 2022), following the approach of the WDC dataset. Having two such datasets could allow for the analysis of both multilingual approaches and generalization of models to other languages. The dataset was derived from popular Polish stores, involving data collection, subsequent cleaning, and transformation into a tabular format compatible with the WDC dataset.

3 Approach & research methodology

In this section, we first name the dataset used in the project, followed by a description of methodologies related to our project. The last subsection contains an explanation of our chosen approach.

To create our solution, we have selected the WDC dataset as the primary dataset for our project. Additionally, we wanted to test our solution on the Polish dataset. However, due to computational and time limitations, we could not utilize this dataset. More so, feature engineering using LLM’s in the experiments on this data only returns English descriptions, automatically translating the Polish text. Therefore, adding these sentences would mix two languages in a single input, which we believed to be a further issue. Thus, the Polish dataset is left as possible future work for this project.

3.1 Related Methodology

Our solution is based on the pipeline in (Tracz et al., 2020). In this article, the authors propose using a transformer architecture, specifically a fine-tuned version of BERT, to embed and then compare a pair of products. As the pre-embedding product representation, a concatenation of the title, attribute values, and units was used. They used a triplet loss objective with the cosine distance for the fine-tuning process. Additionally, various batch construction strategies for selecting the triplets used in training were analyzed.

A similar methodology has been described in (Peeters et al., 2020), where a cross-encoder structure was used instead of a bi-encoder. Additionally, textual representation has been replaced by a concatenation of product brand, title, truncated description and truncated specification table.

An alternative approach was presented in (Peeters and Bizer, 2023), where the similarity score was generated with an LLM under proper prompt construction. However, due to performance limitations present in our problem this solution was disregarded. This lead us to the idea of using LLM’s as a feature engineering tool, which is described further in the report.

Authors of the (Michał Mozdzzonek, 2022) article focused on the Product matching problem and the idea of using transfer learning for data in dif-

ferent languages. Data preparation for the model involved selecting the title column and concatenating it with token markers. This resulted in a single input string for the model, which was then tokenized using a pre-trained model-specific tokenizer.

For the work’s experimental part, using the HuggingFace Transformers library, two types of pre-trained models were used: mBERT and XLM-RoBERTa. The models were pre-trained on Wikipedia articles in about 100 languages. The authors run the models on both WDC and Polish datasets. The F1 score was used as a metric to compare the models.

The mBERT (Devlin et al., 2019) and XLM-RoBERT (Conneau et al., 2020) models consistently outperformed other models. Notably, mBERT excelled, particularly in smaller-sized datasets (small, medium), a trend also observed in the Polish dataset. In the ”Shoes” category, results were slightly lower with the mBERT model. However, XLM-RoBERT performed exceptionally well in ”large” datasets.

These results demonstrate that multilingual models effectively address the product matching problem, often yielding results comparable or superior to prior studies.

3.2 Chosen Approach

Our main aim was to test a method for product similarity matching using reliable language models. An important part of the project is the limitation of the execution time, which is set to 10 ms. Because of that, in our approach, we decided to test out BERT (Bidirectional Encoder Representations from Transformers), DistilBERT (Sanh et al., 2019) as well as RoBERTa (Michał Mozdzonek, 2022). Despite having 40% fewer parameters than the original model, DistilBERT achieves competitive results with the benefit of running 60% faster. RoBERTa improves upon BERT by removing the next sentence prediction objective, utilizing dynamic masking during pre-training, and employing larger mini-batches, enhancing performance in various natural language processing tasks. Additionally, as the project description is not stated otherwise, we assumed that the imposed time limitation doesn’t include feature extraction. Therefore we could utilize Large Language Models (LLMs) as feature extractors. This was important for us because LLMs demonstrated remarkable capabilities

in understanding context capturing complex patterns, and nuances in language.

In the first step we designed a method to evaluate a solution that we were working on. One would want this metric to be able to compare the similarity between many products simultaneously and output a single number, which would measure how good the model is in ranking products on multiple levels. The input vector for the metric consists of cosine similarities between embeddings of product pairs from the transformer models. The resulting metric (called RankedSimilarity) is a sum of Kendall Tau Distance (KDT) (Cicirello, 2019) and Mean Square Error (MSE).

$$\text{RankedSimilarity}(x) = \text{KDT}(x) + \text{MSE}(x, [1, 0.66, 0.33, 0])$$

, where x is a vector of cosine similarity. The former is responsible for keeping similarities of products in order and the latter forces similarities to be equally spaced.

KDT measures the dissimilarity between two rankings by counting the number of pairwise disagreements in the ordering of elements. It considers the number of concordant and discordant pairs (pairs of elements in the same or opposite order in the two rankings) to quantify the similarity or dissimilarity between the rankings. In our case, we ranked cosine similarities of product pairs.

MSE, in our case, measures how much vector of ordered cosine similarities between products pairs are different from the vector of values from 1 to 0 with equal differences between them.

For example: lets say that a, b, c, d are four embeddings of products from transformer model and we know that product b is the most similar to product a and product d is the least similar to product a . Therefore we want the cosine similarities to be in the same order:

$$\cos(a, a) > \cos(a, b) > \cos(a, c) > \cos(a, d)$$

and we represent them as vector:

$$[\cos(a, a), \cos(a, b), \cos(a, c), \cos(a, d)]$$

. KDT assures the order. However using only that, vector of cosine similarities like:

$$[1, 0.99, 0.11, 0.10]$$

would result in perfect score because the values are in order. MSE combats this by forcing a differences of subsequent values to be equal and close as much as possible to:

$$[1, 0.66, 0.33, 0]$$

The main idea behind similarity measure could be shown on a simple example 3.2. As we can see, in all three models the similarities follow expected values - it is higher when compared items are close to each other (like red apple and green apple), and lower when items are different (like red apple and Warsaw University of Technology). All values are high, reaching even 0.9949 for the most different pair of items.

The next step after creating a metric suitable for evaluation was transforming the original dataset to meet our needs. Namely, we extracted a representation of each product using LLM and based on positive pairs, negative pairs and other information about them, we created an evaluation data set in which each observation is list of products arranged in order of similarity.

Prompt engineering is an important technique for getting good results from LLMs (Velásquez-Henao et al., 2023). It involves crafting prompts that explicitly guide the LLM towards the desired output, ensuring that it understands the task and generates the most relevant and informative response. It is an iterative process involving experimentation and refinement to optimize the LLM’s performance for specific tasks. After many iteration we came up with the following prompt: *Given a product title and description, generate a meaningful text representation that captures the essence of the product for effective similarity search. Consider relevant features, attributes, and contextual information to ensure the generated representation reflects the product’s unique characteristics, allowing for accurate comparisons in a similarity search algorithm. Do not answer, just create a representation. TEXT TO REPRESENT: {title+description}*. We drew inspiration from the successful approaches employed in our previous works and we used the prompt along with titles and descriptions of the products to generate augmented text representations of products that will be used as input for transformer models.

Each record in the evaluation dataset consists of representations of 5 products ordered like this $[anchor, anchor', positive, negative, category]$, where:

1. *anchor* is the main product representation that is compared to other products in a record (equivalent of product *a* in earlier example).
2. *anchor'* is a representation of the same product as *anchor* but generated independently from the first one resulting in a different string for the same product
3. *positive* is a text representation of the product from the positive pair from original dataset.
4. *negative* is a text representation of the product from the negative pair from original dataset.
5. *category* is a text representation of a product from different category than original product.

Subsequent products are less and less similar to the *anchor* thanks to which the created dataset is suitable to be used to evaluate the performance of models using the previously described metric.

To train the BERT-like models we utilized Triple Loss that is expressed like this:

$$L(a, b, c) = \max(\|a - b\|_2 - \|a - c\|_2 + \alpha, 0)$$

, where *a* is embedding of an anchor, *b* is embedding of product from positive pair and *c* is an embedding of a product from negative pair. Hyperparameter α is minimal sufficient difference of distances between them.

4 Experiments and Results

In our experiments, we took into account the performance of pre-trained BERT, DistilBERT and RoBERTa models before and after performing our training procedure. We evaluated models using our custom metric described earlier and execution time for inference of each model.

4.1 Exploratory Data Analysis (EDA)

We conducted exploratory data analysis to understand the Web Data Commers dataset better. We focused mainly on basic statistics like number of instances in each class, missing values in each column and length (in words) of text-like features.

Left item	Right item	Similarity level	Cosine Similarity
Red apple	Green apple	Same fruit	0.99230
Red apple	Lemon	Both fruit	0.97707
Red apple	Brick	Both small objects	0.96443
Red apple	Warsaw University of Technology	Not similar	0.93574

Table 1: Example showing intuition behind similarity measuring.

Additionally, we prepared word clouds from titles and descriptions and attempted to divide the pairs but looking into the percentage of words overlap. In data analysis, we mainly focused on features that we planned to extract information from later in the project. Some of the prepared plots are shown in the later part of the section.

A simple plot presented on fig.2 provides information about the observation count in each category and with each label. The label '0' means that the paired products are not similar, label '1' means that the paired products are very similar or the same.

Four plots are presented on one figure fig. 3 for comparison showing the number of missing values in each column in the dataset with additional division into categories (cameras, computers, jewelry, shoes).

To try to distinguish manually (and also out of curiosity) between similar and different products a metric of the percentage of words overlap was proposed. It is calculated by taking two texts and then calculating how much overlap *intersection* and dividing it by the size of a set created from both texts *union*.

Then the distribution of these values was plotted on fig. 4 with differentiating between labels (label=0 - Different products, label=1 - Similar products). This was done for product titles, descriptions and table content.

4.2 Experimental Procedures

In the first step of our experiments we performed comprehensive analysis of the data in order to uncover any key features, trends, and anomalies that may affect our project. It is explained thoroughly in 4.1.

The next step focused on using LLMs to generate augmented text representations. During the project we considered different models as well as services where those models were available. One of the best models is GPT3.5, available via the

OpenAI API. Unfortunately, due to the fact that it is paid, we decided not to use it. We also tried to host the LLama 7B model locally, but the results were not satisfactory. Ultimately, we used LLama 70B provided by HuggingFace for free. However, outputs from the model were in English, regardless of the language in the input, which is one of the reasons we decided to abandon experiments on Polish data set.

In the same time we tried to determine which transformer model would be best suited for the task of calculating products similarity. Due to the project requirements we were forced to consider smaller, more effective models that would provide inference in less than 0.01 seconds per product pair. We decided to use BERT-based models known for their bidirectional architecture and efficient deployment. We planned to utilize HerBERT for the polish dataset, however due to the limitations mentioned earlier we rejected this idea.

Next major step focused on creating evaluation method for multi-level similarity. We recognised two key characteristics that our metrics should have. The first one is ordering of the similar products. It is important from the business perspective, for example placing similar products in a logical order can enhance the overall shopping experience for customers. The second one is to uniformly separate the similarity scores (cosine similarity) to avoid situations where the least similar product has 0.99 similarity score.

After that we came back to LLMs. Using them aims to get richer representation of the product informations. Following previous works we decided to use title and description of the products. We iteratively refined prompt in order to force the model to add some information or descriptions from its knowledge that may be helpful in further processing.

With the Hugging Face API, a prompt, and a metric to evaluate, we prepared a dataset that fits our needs. Each row consisted of five text representations of the products generated by LLM.

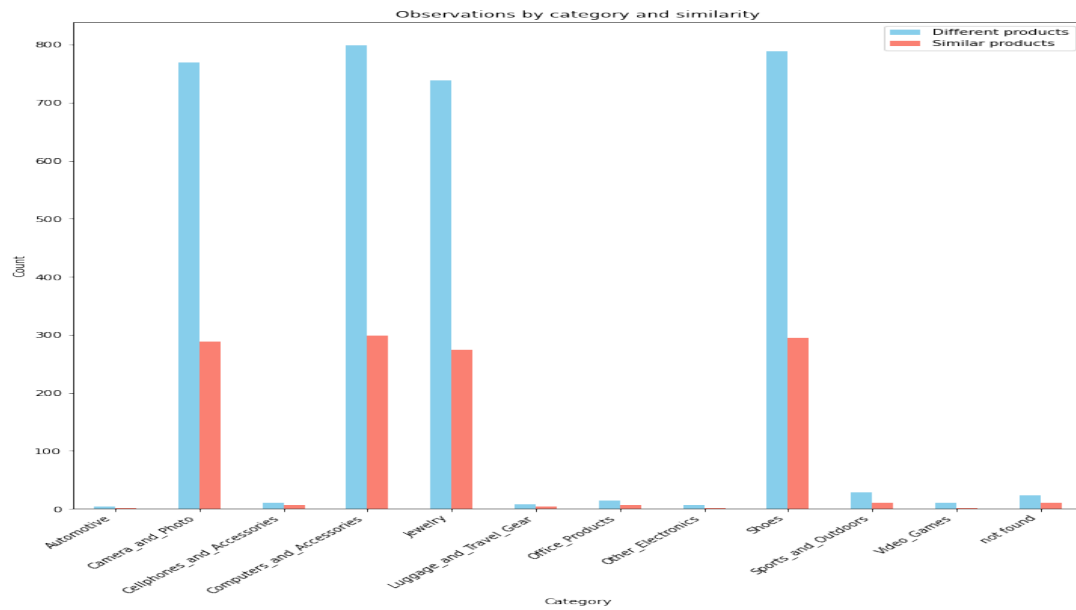


Figure 2: Number of o observations in each category and label

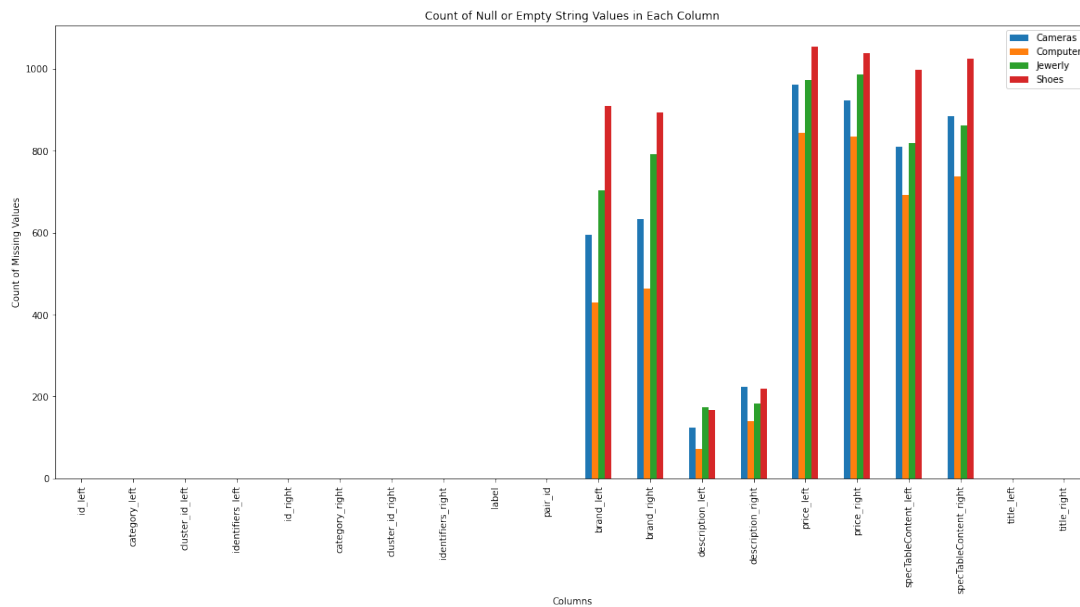


Figure 3: Number of missing values in each column in the dataset with additional division into categories (cameras, computers, jewelry, shoes).

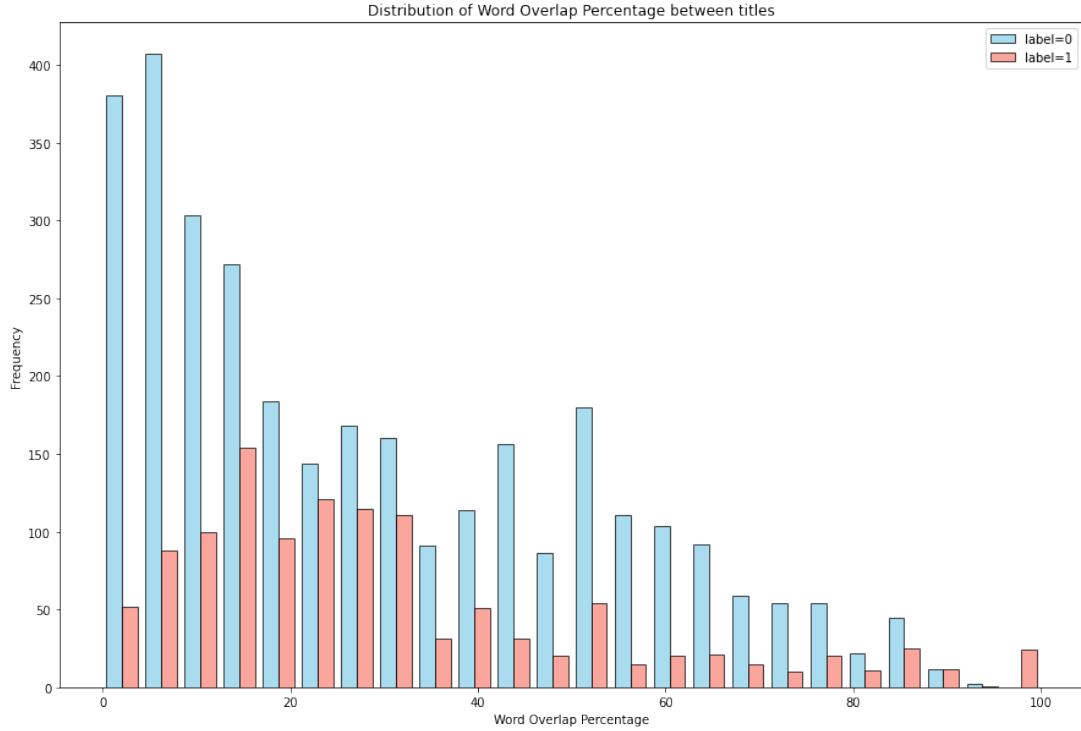


Figure 4: Percentage of words overlap in titles.

Each column represented different level of similarity to the "main" product. For more information look 3.2.

Using our augmented dataset we performed fine-tuning to leverage the knowledge gained during the initial training.

Our final results consist of 5 different BERT models performance on the WDC dataset. As seen on figure 5 all tested models achieve high results of the measure, however it is worth noting that the finetuned models score worse than the pretrained ones. This issue is further described in section ??.

The best model (BERT) is closely followed by the second best (DistilBERT), however a significant difference is measured in the inference time of the models, as seen on 5. The distilled version executes nearly twice as fast, while barely lagging behind in the metric score.

5 Discussion on your results

Our results are far from perfect, however they're satisfactory as the first attempt at defining the measure for incorporating taxonomy into a similarity comparison. Adding LLM data augmentation in this manner is a significant, novel approach (as far as we know) to elevate the embeddings accuracy. The worse performance of

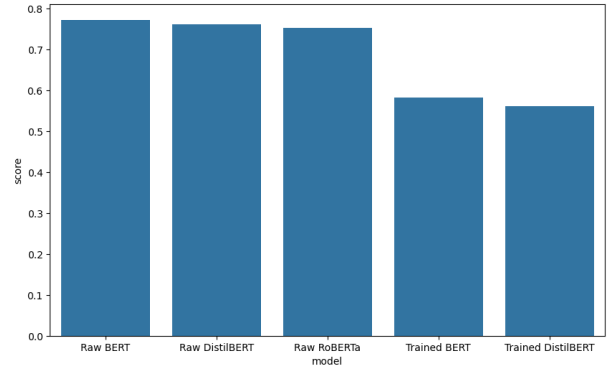


Figure 5: Metric score for pre-trained models before and after our training

fine tuned models is a surprising effect, which might be caused by the low amount of data available in the dataset or the big concept gap between tasks of similarity measurement and prediction of masked sentence elements (as in BERT). The untrained models, on the same example, achieved an example cosine similarity vector like $([1.0000, 0.9597, 0.9527, 0.9481])$ while the fine-tuned model $([0.9417, 0.9588, 0.9497, 0.9528])$. This shows how exactly these model fail since with training it approach the ideal score of $([1, 0.66, 0.33, 0])$ rather than get further away from it and lose the decreasing characteristic.

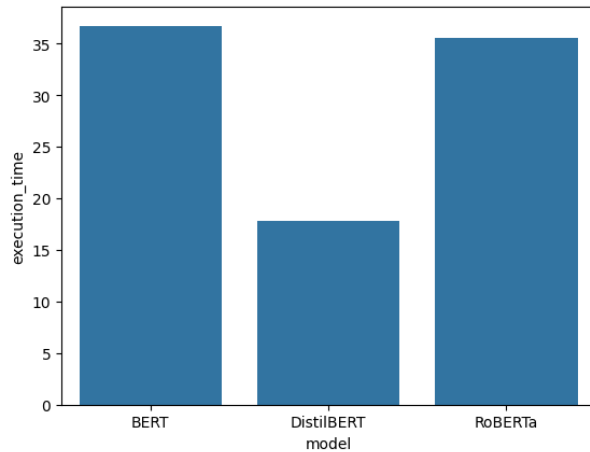


Figure 6: Inference time for chosen models.

While we set out to hopefully reach the 10ms per comparison threshold, it is not a significant failure. The fastest-performing model scored similarly to the slowest but highest metric model, achieving 18ms per pair. This time consists of both the embedding of text and similarity measurement, however the reason why it isn't that much of an issue is the different approach taken in standard e-commerce projects. The data for a single product would be embedded once and saved in an indexed vector database. This would allow for fast retrieval of the data when needed, and the time of similarity measurement is irrelevant, as it reaches $46.6\mu s \pm 4.64\mu s$ (on 700000 tests) which is around 400 times less. Additional ways of decreasing the inference time are further described in future works ??.

Additionally, in our research we encountered multiple issues with the limited computational resources, like:

- Difficulties with setting up the dataset nor the LLM on EDEN (problems with the available RAM or drivers or time, EDEN had problems for a significant amount of time).
- Inference time of self-hosted LLM turned out to be extremely long.
- Inference time of publicly hosted LLM also turned out to be very long.
- API limitations ran out fast.
- Free GPU for training run out fast.

6 Conclusions and future work

In future developments of our e-commerce recommendation systems project, several key areas exist to explore. Firstly, expanding our methodology to include a broader range of datasets is crucial for assessing the generalizability of our product similarity matching system across different domains. Testing the solution on a bigger and cleaner dataset could provide the necessary amount data for the models, as scaling typically improves the results of neural networks.. Additionally, testing our recommendation system on advanced Large Language Models (LLMs) beyond the current state-of-the-art models like GPT-3 or BERT will ensure adaptability to evolving NLP technologies. Experimenting with various prompts for feature extraction and evaluating different system prompts during the recommendation process will contribute to optimizing the system's performance. Furthermore, conducting comparative analyses without the reliance on Large Language Models will provide insights into the specific impact of these models on the recommendation system. Exploring alternative methods for description extraction, including domain-specific techniques or specialized models, is essential for diversifying the approach. Finally, testing the performance and success of the recommendation system in a real e-commerce environment, potentially through collaboration with industry partners or deployment in controlled settings, will validate its practical applicability and reveal opportunities for refinement and enhancement. Continuing research in these areas will contribute to the evolution of e-commerce recommendation systems, ensuring improved user experiences and advancements in the field.

References

- [Bobadilla et al.2013] Jesús Bobadilla, Fernando Ortega, Antonio Hernando, and Abraham Gutiérrez. 2013. Recommender systems survey. *Knowledge-based systems*, 46:109–132.
- [Brown et al.2020] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901.
- [Cicirello2019] Vincent A Cicirello. 2019. Kendall tau sequence distance: Extending kendall

Task (Total time est.)	Main Contributor	Secondary Contributor
Project Proposal (10h)	Everyone	-
Exploratory Data Analysis (6h)	<i>Kinga Ulasik</i>	<i>Szymon Rećko</i>
Proof Of Concept (6h)	<i>Mateusz Sperkowski</i>	<i>Patryk Tomaszewski</i>
First Milestone Presentation & Raport (14h)	Everyone	-
Reviews 1 (4h)	Everyone	-
Generating the additional descriptions (4h)	<i>Szymon Rećko</i>	<i>Mateusz Sperkowski</i>
Tuning the models (8h)	<i>Patryk Tomaszewski</i>	<i>Kinga Ulasik</i>
Reviews 2 (4h)	Everyone	-
Second Milestone Presentation & Raport (19h)	Everyone	-

Table 2: Work contribution and estimated total time for each task. While two main contributors are listed, all tasks were done with the support and help of everyone in the team when necessary.

- tau from ranks to sequences. *arXiv preprint arXiv:1905.02752*.
- [Conneau et al.2020] Alexis Conneau, Kartikay Khan-delwal, Naman Goyal, Vishrav Chaudhary, Guil-laume Wenzek, Francisco Guzmán, Édouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2020. Unsupervised cross-lingual representation learning at scale. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 8440–8451.
- [Devlin et al.2019] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *North American Chapter of the Association for Computational Linguistics*.
- [Gou et al.2021] Jianping Gou, Baosheng Yu, Stephen J Maybank, and Dacheng Tao. 2021. Knowledge dis-tillation: A survey. *International Journal of Com-puter Vision*, 129:1789–1819.
- [Kasneci et al.2023] Enkelejda Kasneci, Kathrin Seßler, Stefan Küchemann, Maria Bannert, Daryna Dementieva, Frank Fischer, Urs Gasser, Georg Groh, Stephan Günnemann, Eyke Hüllermeier, et al. 2023. Chatgpt for good? on opportunities and challenges of large language models for education. *Learning and individual differences*, 103:102274.
- [Koch et al.2015] Gregory Koch, Richard Zemel, Rus-lan Salakhutdinov, et al. 2015. Siamese neural networks for one-shot image recognition. In *ICML deep learning workshop*, volume 2. Lille.
- [Liu et al.2018] Zhuang Liu, Mingjie Sun, Tinghui Zhou, Gao Huang, and Trevor Darrell. 2018. Re-thinking the value of network pruning. In *International Conference on Learning Representations*.
- [Michał Mozdzonek2022] Sergiy Tkachuk Szymon Łukasik Michał Mozdzonek, Anna Wróblewska. 2022. Multilingual trans-formers for product matching – experiments and a new benchmark in polish.
- [Peeters and Bizer2023] Ralph Peeters and Christian Bizer. 2023. Using chatgpt for entity matching. *arXiv preprint arXiv:2305.03423*.
- [Peeters et al.2020] Ralph Peeters, Christian Bizer, and Goran Glavaš. 2020. Intermediate training of bert for product matching. *small*, 745(722):2–112.
- [Polino et al.2018] Antonio Polino, Razvan Pascanu, and Dan-Adrian Alistarh. 2018. Model compres-sion via distillation and quantization. In *6th Inter-national Conference on Learning Representations*.
- [Sanh et al.2019] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. 2019. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*.
- [Scarselli et al.2008] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini. 2008. The graph neural net-work model. *IEEE transactions on neural networks*, 20(1):61–80.
- [Thirunavukarasu et al.2023] Arun James Thirunavukarasu, Darren Shu Jeng Ting, Kabi-lan Elangovan, Laura Gutierrez, Ting Fang Tan, and Daniel Shu Wei Ting. 2023. Large language models in medicine. *Nature medicine*, 29(8):1930–1940.
- [Tracz et al.2020] Janusz Tracz, Piotr Iwo Wójcik, Kalina Jasinska-Kobus, Riccardo Belluzzo, Robert Mroczkowski, and Ireneusz Gawlik. 2020. Bert-based similarity learning for product matching. In *Proceedings of Workshop on Natural Language Pro-cessing in E-Commerce*, pages 66–75.
- [Vaswani et al.2017] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural infor-mation processing systems*, 30.
- [Velásquez-Henao et al.2023] Juan David Velásquez-Henao, Carlos Jaime Franco-Cardona, and Lorena Cadavid-Higuita. 2023. Prompt engineering: a

methodology for optimizing interactions with ai-language models in the field of engineering. *DYNA*, 90(230):9–17.

[Wei et al.2022] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. 2022. Emergent abilities of large language models. *Transactions on Machine Learning Research*.

[Wu et al.2023] Shijie Wu, Ozan Irsoy, Steven Lu, Vadim Dabravolski, Mark Dredze, Sebastian Gehrmann, Prabhanjan Kambadur, David Rosenberg, and Gideon Mann. 2023. Bloomberggpt: A large language model for finance. *arXiv preprint arXiv:2303.17564*.

[Zhang et al.2019] Shuai Zhang, Lina Yao, Aixin Sun, and Yi Tay. 2019. Deep learning based recommender system: A survey and new perspectives. *ACM computing surveys (CSUR)*, 52(1):1–38.