

JointMap: Joint Query Intent Understanding For Modeling Intent Hierarchies in E-commerce Search

Ali Ahmadvand*
Emory University
ali.ahmadvand@emory.edu

Faizan Javed
The Home Depot
faizan_javed@homedepot.com

Surya Kallumadi
The Home Depot
surya@ksu.edu

Eugene Agichtein
Emory University
eugene.agichtein@emory.edu

Abstract

An accurate understanding of a user’s query intent can help improve the performance of downstream tasks such as query scoping and ranking. In the e-commerce domain, recent work in query understanding focuses on the query to product-category mapping. But, a small yet significant percentage of queries (in our website 1.5% or 33M queries in 2019) have non-commercial intent associated with them. These intents are usually associated with non-commercial information seeking needs such as discounts, store hours, installation guides, etc. In this paper, we introduce Joint Query Intent Understanding (JointMap), a deep learning model to simultaneously learn two different high-level user intent tasks: 1) identifying a query’s commercial vs. non-commercial intent, and 2) associating a set of relevant product categories in taxonomy to a product query. JointMap model works by leveraging the transfer bias that exists between these two related tasks through a joint-learning process. As curating a labeled data set for these tasks can be expensive and time-consuming, we propose a distant supervision approach in conjunction with an active learning model to generate high-quality training data sets. To demonstrate the effectiveness of JointMap, we use search queries collected from a large commercial website. Our results show that JointMap significantly improves both “commercial vs. non-commercial” intent prediction and product category mapping by 2.3% and 10% on average over state-of-the-art deep learning methods. Our findings suggest a promising direction to model the intent hierarchies in an e-commerce search engine.

ACM Reference Format:

Ali Ahmadvand, Surya Kallumadi, Faizan Javed, and Eugene Agichtein. 2020. JointMap: Joint Query Intent Understanding For Modeling Intent Hierarchies in E-commerce Search. In *Proceedings of the 43rd International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR ’20)*, July 25–30, 2020, Virtual Event, China. ACM, Taipei, Taiwan, 4 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

*Research was conducted while interning at The Home Depot Search & NLP team.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

SIGIR ’20, July 25–30, 2020, Virtual Event, China

© 2020 Association for Computing Machinery.
ACM ISBN 978-x-xxxx-xxxx-x/YY/MM... \$15.00
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

| Search Queries | intent | Product Categories |
|-------------------------------------|----------------|-------------------------------|
| where is my shipped order | non-commercial | - |
| how to install my tiles | non-commercial | - |
| cost to rent a carpet cleaner | non-commercial | - |
| 18 volt ryobi | commercial | [tools, electrical, lighting] |
| 24 in. classic Samsung refrigerator | commercial | [appliance, electrical] |

Table 1: Dataset sample queries and their associated labels.

1 INTRODUCTION AND RELATED WORK

Query intent understanding is a key step in designing advanced retrieval systems like e-commerce search engines [6]. Various approaches have been proposed to address query understanding such as 1) considering predefined high-level categories (i.e., informational, navigational, and transactional), 2) deploying semi-supervised learning with click graphs, 3) considering temporal query intent modeling, 4) understanding word-level user intent, and 5) applying relevance feedback and user behaviors. Although there has been a significant improvement in user intent inference, query understanding remains a major challenge [13].

E-commerce search queries have multiple intents associated with them. Ashkan et al. [1] categorized search queries for e-commerce websites into commercial and non-commercial intents. However, Zhao et al. [14] ignore the non-commercial queries due to small percentage of the search traffic. Commercial queries are queries with purchasing intent, while non-commercial queries cover a wide range of customer services (e.g., “military discounts” and “installation guides”) as shown in Table 1.

Query understanding in e-commerce search is challenging: 1) queries are often short, vague, and suffer from the lack of textual evidence [7], 2) small variation in textual evidence causes a drastic change in query intent; for example, “30 in. 5.8 cu. ft. gas range installation kit” has commercial intent but “30 in. 5.8 cu. ft. gas range installation”, has non-commercial intent, 3) product category mapping is a multi-label and non-exclusive problem. A practical solution must include a broader possible set of correct categories, while simultaneously keeping precision as high as possible [14], 4) there is class imbalance in both commercial vs. non-commercial and product category mapping tasks, because only a small fraction of data (1.5% in our domain) has a non-commercial intent, and within the commercial queries, some product categories contain significantly more samples compared to others, and 5) commercial queries are easy to identify using user behavior information like click rates; however, that is not the case for non-commercial queries.

To address these problems, we introduce a new method of jointly learning query intent and category mapping, which allows transferring the inductive bias between these two relevant tasks. Also,

we leverage label representation, which provides a richer representation to model the product categories. Finally, we propose an active learning algorithm to generate data for commercial vs. non-commercial intent. To address the class imbalance problem, we deploy focal loss, which is borrowed from computer vision.

Joint learning has been proposed as a practical approach to simultaneously learn relevant tasks due to the transfer of the inductive bias among them. Joint-learning finds applications in computer vision and natural language understanding [8]. Joint-learning improves the regularization and generalization of the learning models by utilizing the domain information [3]. In addition, with a joint model that addresses multiple tasks, only one model needs to be deployed; this contributes to reducing overhead and facilitates the maintenance of the system [12]. In this paper, we propose a joint-learning model that simultaneously learns both commercial and non-commercial query intents, and maps the incoming commercial query to a set of relevant product categories.

In this paper, we introduce a data-driven approach, which we call Joint Query Intent Mapping (JointMap). JointMap leverages the label representation proposed by Guoyin et al. [11] and modifies it to be applicable for a joint-learning task. JointMap also utilizes self-attention mechanism to improve the quality of the joint word-label attention vectors. For product category mapping, JointMap handles the imbalanced class problem using focal loss [9] which has been well-studied in the computer vision field to control the sparse set of candidate object locations. Finally, we propose an approach based on distant supervision in combination of active learning to generate both commercial and non-commercial queries.

In summary, our contributions are: 1) proposing a deep learning model to jointly learn product category mapping as well as users’ non-commercial intents, 2) developing an active learning algorithm in conjunction with distant supervision to generate a user intent dataset from e-commerce data logs, and 3) modifying the joint word-category representation for query intent mapping tasks in e-commerce, as described in detail next.

2 MODEL OVERVIEW

In this section, we present the network architecture of JointMap, as shown in Figure 1. JointMap utilizes both word and category embeddings in which both representations are jointly trained to achieve an efficient semantic representation for a query. The proposed model consists of two deep learning layers: the first layer for the understanding of the user’s commercial intent and the second layer for the prediction of relevant product categories in the taxonomy. As a result, the proposed model contains three embedding layers: a word embedding layer and two category embeddings layers, i.e., commercial vs. non-commercial and product-categories. Both category embedding types are concatenated, to compute the final product category representations. Then, a Compatibility Matrix (CM) is generated by computing the cosine similarity between the label and word representations. CM represents the relative spatial information among consecutive words (phrases) with their associated product category and commercial vs. non-commercial labels. Finally, CM is passed through a Multi-head self-attention layer to calculate attention scores. The word vectors simultaneously go through two Highway layers, and the output of each Highway is

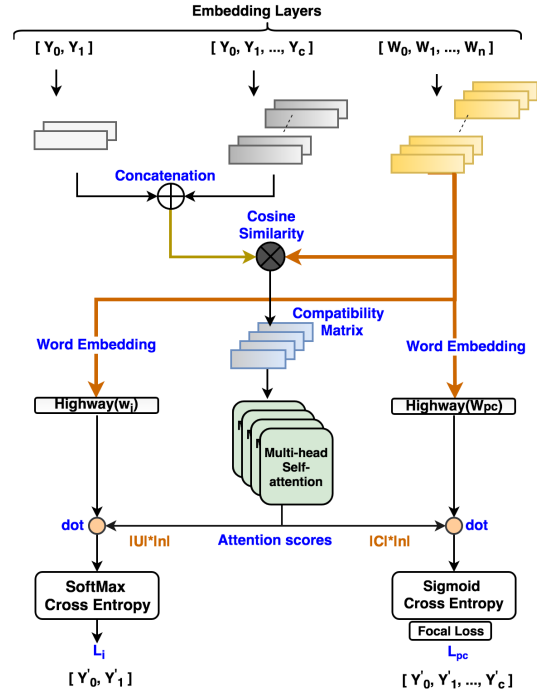


Figure 1: JointMap network architecture.

multiplied by their corresponding attention scores to generate the final query representation. Finally, the loss value of \mathcal{L}_{pc} is computed using sigmoid cross-entropy for the product category mapping. Also, the loss value \mathcal{L}_i is calculated using Softmax cross-entropy for determining the query’s commercial intent.

In the next section we explain the details of the proposed model.

2.1 Joint-Learning of High-Level Intent Tasks

We now introduce JointMap, a joint-learning model for high-level user intent prediction.

Suppose there is a search query dataset $D = \{Q, C, U\}$, where Q is a set of search queries, U represents user commercial vs. non-commercial intent, and C is the candidate product category set. Each query consists of a sequence of words $q = [w_1; w_2; \dots; w_n]$ of size $n = 10$, and represents as $\mathbf{W}^{|W| \times V}$. Also, C and U are mapped to the embedding spaces $\mathbf{C}^{|C| \times V}$ and $\mathbf{U}^{|U| \times V}$, respectively. Then, the matrices C and U are concatenated to illustrate the whole label space. The word and label embeddings are initialized with Word2Vec and random embeddings of size $|V| = 300$, respectively. Cosine similarity between L and W is computed for each query q to extract the relative spatial information among consecutive words with their associated labels, where \otimes indicates the cosine similarity function.

$$\mathbf{H} = (\mathbf{C} + \mathbf{U})^{(|C|+|U|) \times V} \otimes (\mathbf{W}^{n \times V})^T \quad (1)$$

To extract the contribution of the words concerning their category, a multi-head self-attention mechanism with n different heads is implemented on H . Multi-head self-attention contains a parallel of linear projections of a single scaled dot-product function. Eq. 2 shows a single head of the self-attention mechanism.

$$\mathbf{G} = \text{SoftMax}\left(\frac{\mathbf{H}\mathbf{K}^T}{\sqrt{d_k}}\right) \mathbf{V} \quad (2)$$

where K is the key matrix, V is the value matrix, and d_k is the dimension of the keys. Also, each projection is responsible for extracting the attention between word-label in a query and computes using weighted sum of the values. Next, G is split into two matrices of size $\hat{G} = (|C| \times n)$ and $\widehat{G} = (|U| \times n)$. For both tasks, the word embedding vectors W are fed into a highway encoder layer, which has shown its effectiveness in improving network capacity [13]. Then, the output is multiplied by their corresponding attention scores of \hat{G} .

$$\alpha_1 = Highway_1(W), \alpha_2 = Highway_2(W) \quad (3)$$

$$\alpha_i = \text{sigmoid}(r_w) \rightarrow r_w = \text{relu}(\text{word2vec}(w)) \quad (4)$$

$$W_1 = \sum_{i=1}^n \hat{G}_i \times \alpha_1, W_2 = \sum_{i=1}^n \widehat{G}_i \times \alpha_2 \quad (5)$$

Then, resulted W_1 and W_2 have the size of $(n \times V)$. They go through a fully connected layer to generate the semantic representations of both tasks. For product category mapping, a sigmoid cross-entropy loss function \mathcal{L}_{pc} is used since in sigmoid, the loss computed for every output s_i is not affected by other component values. Also, a binary softmax cross-entropy loss \mathcal{L}_i is applied to train the user commercial vs. non-commercial intent.

$$\mathcal{L}_{pc} = - \sum_{c=1}^{|C|} t_c \log(\text{Sigmoid}(s_c)) \quad (6)$$

$$\mathcal{L}_i = -t_1 \log(\text{SoftMax}(s_1)) - (1 - t_1) \log(1 - \text{SoftMax}(s_1)) \quad (7)$$

Where s_c represents the prediction distribution and t_c indicates the target labels. To address the class imbalance problem, particularly in the product category dataset, we update the loss values based on focal loss proposed in [9]. The focal loss was initially proposed for object detection and removing the effect of extreme foreground-background class imbalance in the images.

$$\mathcal{L}_{focal_{pc}} = \sum_{c=1}^{|C|} \alpha_c (\text{Sigmoid}(s_c) - t_c)^Y \log(\text{Sigmoid}(s_c)) \quad (8)$$

where t is the target vector, c is the class index, and $(f(s) - t)^Y$ is a factor to decrease the influence of well-classified samples.

JointMap overall loss: The final loss function is computed using a weighted loss over commercial vs. non-commercial, product category mapping intents.

$$\mathcal{L}_{total} = \beta_1 \mathcal{L}_{focal_{pc}} + \beta_2 \mathcal{L}_i \quad (9)$$

3 DATASET OVERVIEW

In this section, we describe the dataset collected from search logs of a large e-commerce search engine in July 2019, and provide details the algorithms used for generating user-intent datasets. We propose an algorithm to simultaneously generate both datasets, which consists of three steps: 1) generating the commercial vs. non-commercial queries, 2) oversampling of the non-commercial queries to balance the dataset, and 3) creating the product category dataset based on the commercial queries. Algorithm. 1 represents the steps for generating commercial vs. non-commercial samples. In this method, first we need to generate a small-size dataset that covers all expected non-commercial intents (e.g., "installation guides").

Then, we over-sample the non-commercial queries as described in [4] to make the dataset balanced (only 1.5% of the queries have a non-commercial intent). Similar to [14], we utilize user behavior

Result: *Commercial Vs. Non-commercial Dataset*

D_init = A small-size Dataset by human supervision ;

Test = Hold-out test set;

while Accuracy < threshold **do**

 D = Expand(D_init) using KNN;

 Confidence Scores = SVM(D);

 TS = Find(tricky samples) using confidence scores;

 D = Re-label(TS) using human supervision;

 D_init = D;

 Accuracy = Compute_Accuracy(Test);

end

Algorithm 1: Commercial vs. non-commercial dataset.

data like click rate, to generate the category labels associated with each commercial query. Algorithm. 2 describes different steps to create the product mapping dataset.

Result: *Product Category Dataset*

Product_category = {};

for each query in Q **do**

 pid_list = Extract(pid that user clicks)

for pid in pid_list **do**

 category_list = Find(category(pid) in taxonomy)

end

for category in category_list **do**

if if click_rate > r **then**

 product_category(query).add(category)

end

end

end

Algorithm 2: Product category dataset generator.

Finally, a dataset of size 195K with 32 product categories such as *tools, appliance, outdoors, etc.* extracted from the search logs.

4 EXPERIMENTAL SETUP

In this section, we describe the parameter setting, metrics, baseline models, and experimental procedures used to evaluate JointMap.

Parameter Setting. We used Adam optimizer with a learning rate of $\eta = 0.001$ and a mini-batch of size 64 for training. The dropout rate of 0.5 is applied at the fully-connected and ReLU layers to prevent the model from overfitting.

Evaluation Metrics. To evaluate JointMap, both Micro- and Macro-averaged F1-score for both tasks are reported.

Methods Compared. We summarize the **multi-label** classification methods compared in the experimental results.

- **Tf*idf + SVM:** One-Vs-Rest SVM with a linear kernel.
- **VDCNN:** Very Deep Convolutional Neural Network [5].
- **FastText:** Text classification method developed by Facebook[2].
- **LEAM:** Word-label representation model[11].
- **XML-CNN:** Extreme multi-label text classification [10].
- **JointMap:** The proposed model.

Dataset Experimental Design. We use an SVM model with n-gram tf*idf as features to perform distant supervision method due to multiple reasons: 1) SVM is fast and scalable, 2) the features and results are interpretable for supervisors, 3) SVM has proved

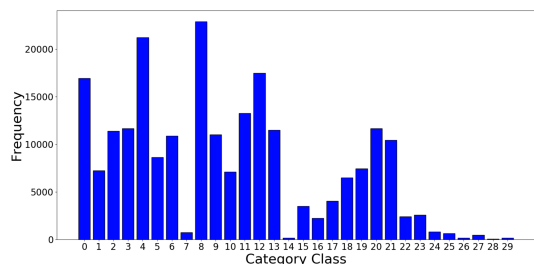


Figure 2: Product category distribution.

its effectiveness on text data, 4) SVM provides confidence scores to detect the tricky samples. Moreover, two different human annotators were asked to label 540 samples manually. The (Matching, Kappa) scores of (0.98, 0.96) are computed, which is a significant agreement. The category distribution is shown in Figure 2.

4.1 Main Results and Ablation Analysis

To evaluate the models described in Section 4, 70% of the dataset is used for training, 10% for validation, and 20% for test. Table 2 summarizes the performance of the models. The results are reported for both commercial vs. non-commercial classification and product category mapping. All the improvements are statistically significant using a one-tailed Student’s t-test with a p-value < 0.05.

| Method | Dataset | | | |
|--------------|-------------------------------|---------------|--------------------------|--------------|
| | Commercial vs. Non-commercial | | Product Category Mapping | |
| | Macro-F1 | Micro-F1 | Macro-F1 | Micro-F1 |
| tf*idf+SVM | 90.71 | 90.26 | 48.75 | 76.84 |
| VDCNN [5] | 91.28 | 91.34 | 51.41 | 79.34 |
| FastText [2] | 92.18 | 92.15 | 60.06 | 79.69 |
| XML-CNN [10] | 93.11 | 93.01 | 58.40 | 81.61 |
| LEAM [11] | 93.96 | 93.66 | 58.90 | 81.31 |
| JointMap | 94.80 (+1.1%) | 94.63 (+1.0%) | 62.60 (+6.3%) | 83.01 (2.1%) |

Table 2: Macro- and Micro- averaged F1 for different models. The improvements reported against LEAM.

For the user commercial intent mapping task, the results indicate that the Macro-averaged F1 improves 4.5%, 3.8%, 2.8%, 1.0%, and 1.8% compared to tf*idf, VDCNN, FastText, LEAM, and XML-CNN models respectively. In product category mapping task, the improvements are more significant. There is improvement of 28.4%, 22.1%, 4.2%, 6.3%, and 7.2% over tf*idf, VDCNN, FastText, LEAM, and XML-CNN models, respectively. As a result, JointMap improves macro-averaged F1 scores over current state-of-the-art deep learning models by 2.3% on commercial vs. non-commercial intents, and a 10% improvement over product category mapping.

In reference to user commercial intent prediction, a 2.3% improvement is considerable since it is in the context of a large e-commerce search engine that receives billions of search queries per year. For product category mapping, the F1-averaged macro experiences a higher jump when compared to the F1-averaged micro (6.3% vs. 2.1%). This improvement indicates the positive impact of inductive bias between these two tasks, which not only boosts the performance of majority classes, but it also contributes to minority classes. For instance, the Macro-average F1 for 8-button minority classes shows in Figure. 2 for XML-CNN and LEAM are 21.76% and 18.33%, respectively, while this number jumps to 31.28% for JointMap.

Focal Loss Impact. Using focal loss deteriorates the overall micro- and macro- averaged F1-scores by 0.6%, 1.5%, respectively. However, the macro-average F1 on 8-button minority classes without focal loss is 31.28%, while with presence of focal loss is 33.81%. This shows a relevant improvement of 8.1%. Also, we observe that in absence of focal loss, the performance of at least two of the minority classes is 0%, therefore making the use of focal loss necessary.

Parameter Tuning. To evaluate the impact of hyper-parameter tuning in JointMap, we implemented a grid search approach on β_1 and β_2 in Eq. 9. We observed that using a smaller β for each task causes a slower convergence for that specific task. However, the final results is not significantly different. In our experiments, a simple average works as good as a fine-tuned hyper-parameter model. For focal loss hyper-parameter tuning, we repeat the experiments with different γ values of 1, 1.2, 1.5, and 2. We observed that the best results achieve using the $\gamma = 1.5$, where the original paper suggested using $\gamma = 2$ for computer vision application.

5 CONCLUSIONS AND FUTURE WORK

We introduced JointMap, a deep learning model designed for jointly learning two high-level intent tasks on e-commerce search data. JointMap utilized word and label representations and leveraged focal loss to tackle class imbalance problem in catalog categories. Our results were promising compared to the state-of-the-art deep learning models with an average raise of 2.3% and 10.9% on Macro-averaged F1 in user commercial vs. non-commercial intent and product category mapping, respectively. Our future work includes tuning the JointMap model incorporate contextual information within a session. In summary, the presented work advances the state-of-the-art user intent prediction, and lays the groundwork for future research on user intent understanding in e-commerce.

Acknowledgements We gratefully acknowledge the financial and computing support from The Home Depot Search & NLP team.

References

- [1] A. Ashkan, C. L. Clarke, E. Agichtein, and Q. Guo. Classifying and characterizing query intent. In *proceedings of ECIR*, pages 578–586. Springer, 2009.
- [2] P. Bojanowski, E. Grave, A. Joulin, and T. Mikolov. Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146, 2017.
- [3] R. Caruana. Multitask learning. *Machine learning*, 28(1):41–75, 1997.
- [4] F. Charte and et al. Dealing with difficult minority labels in imbalanced multilabel data sets. *Neurocomputing*, 326:39–53, 2019.
- [5] A. Conneau, H. Schwenk, L. Barrault, and Y. Lecun. Very deep convolutional networks for text classification. *arXiv preprint arXiv:1606.01781*, 2016.
- [6] W. B. Croft, M. Bendersky, H. Li, and G. Xu. Query representation and understanding workshop. In *SIGIR Forum*, volume 44, pages 48–53, 2010.
- [7] J.-W. Ha, H. Pyo, and J. Kim. Large-scale item categorization in e-commerce using multiple recurrent neural networks. In *SIGKDD*, pages 107–115. ACM, 2016.
- [8] C. Khatri, R. Goel, B. Hedayatnia, A. Metanillou, A. Venkatesh, R. Gabriel, and A. Mandal. Contextual topic modeling for dialog systems. In *2018 IEEE Spoken Language Technology Workshop (SLT)*, pages 892–899. IEEE, 2018.
- [9] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *ICCV*, pages 2980–2988, 2017.
- [10] J. Liu, W.-C. Chang, Y. Wu, and Y. Yang. Deep learning for extreme multi-label text classification. In *proceedings of SIGIR*, pages 115–124, 2017.
- [11] G. Wang, C. Li, W. Wang, Y. Zhang, D. Shen, X. Zhang, R. Henao, and L. Carin. Joint embedding of words and labels for text classification. *arXiv preprint arXiv:1805.04174*, 2018.
- [12] J. Wang, J. Tian, and et al. A multi-task learning approach for improving product title compression with user search log data. In *proceeding of AAAI*, 2018.
- [13] H. Zhang and et al. Generic intent representation in web search. In *proceedings of SIGIR*, pages 65–74. ACM, 2019.
- [14] J. Zhao, H. Chen, and D. Yin. A dynamic product-aware learning model for e-commerce query intent understanding. In *CIKM*, pages 1843–1852, 2019.