

Taxonomic Recommendations of Real Estate Properties with Textual Attribute Information

ZACHARY HARRISON* and ANISH KHAZANE*, Zillow Group, USA

In this extended abstract, we present an end to end approach for building a taxonomy of home attribute terms that enables hierarchical recommendations of real estate properties. We cover the methodology for building a real-estate taxonomy, metrics for measuring this structure’s quality, and then conclude with a production use-case of making recommendations from search keywords at different levels of topical similarity.

CCS Concepts: • **Information systems** → **Ontologies; Personalization; Language models; Retrieval models and ranking; Recommender systems.**

Additional Key Words and Phrases: taxonomy, recommender systems, data mining

ACM Reference Format:

Zachary Harrison and Anish Khazane. 2022. Taxonomic Recommendations of Real Estate Properties with Textual Attribute Information. In *Sixteenth ACM Conference on Recommender Systems (RecSys '22)*, September 18–23, 2022, Seattle, WA, USA. ACM, New York, NY, USA, 8 pages. <https://doi.org/10.1145/3523227.3547386>

1 INTRODUCTION AND MOTIVATION

One of Zillow’s primary objectives is to provide customers with easily accessible and relevant information for millions of listings on our platform. Our recently launched home insights feature enables this vision by highlighting the most unique keyphrases from a property’s description on their respective home details page.

These tags can cover a wide array of attribute types, such as interior (e.g hardwood floors), exterior (e.g outdoor patio) or even community or location attributes (e.g nearby golf course) that can be useful for users searching for specific features in homes. Categorizing properties simply based on these raw phrases is difficult due to the significant variety in possible topics. However, organizing these terms in a taxonomy structure can not only help with categorizing these terms under different levels of topical similarity but also enable hierarchical recommendations.

In this extended abstract, we present an approach for constructing a real-estate specific taxonomy of home attribute terms that can enable hierarchical recommendations of property listings. We first describe a semi-automated approach to construct this taxonomy, propose metrics for measuring edge quality, then conclude with a production use-case of making hierarchical recommendations from search keywords at different levels of topical similarity. We hope this work can contribute to the few examples of tree-based recommender systems used in industry such as Pinterest’s pin2interest knowledge graph system and Yahoo’s taxonomy powered recommender systems for modeling user purchase behavior [2, 4].

*Both authors contributed equally to this research.

Permission to make digital or hard copies of part or all 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 third-party components of this work must be honored. For all other uses, contact the owner/author(s).

© 2022 Copyright held by the owner/author(s).

Manuscript submitted to ACM

Table 1. State of taxonomy after each round of expansion detailed in Section 2.1-2.2. Link pruning and human-in-the-loop revision is critical for constructing a deeper taxonomy structure for hierarchical recommendations.

Expansion Stage	# Nodes	# Edges	# Parents	# Leaf Nodes	Max Depth
Seed taxonomy	2560	2559	50	2509	2
Embedding Clustering	4023	4022	50	3972	2
Link Pruning and Manual Revision	9138	9137	806	8332	7

2 METHODOLOGY AND CHALLENGES

2.1 Bootstrapping a Seed Taxonomy

We begin by using a pretrained scene entity detection network to label thousands of photos spanning many types of properties (e.g single family homes, apartment rentals, condos) with high-level concept attributes like “KITCHEN” or “BATHROOM”, and also a generative model to tag images with keywords (e.g “granite countertops”, “wood flooring”) describing the contents of the image [5]. Thus, we can rapidly bootstrap a 2-level taxonomy by mapping the keywords of an image to the high-level scene entities identified by the deep network. This approach yields 2560 nodes and 2559 edges as displayed in Table 1.

We then expand this structure by training a fasttext embedding model on listing descriptions to generate subword embedding representations of the aforementioned home insight phrases and scene entities in the taxonomy [1]. Following embedding generation, we initialize a k-nearest neighbors model with the scene entity embeddings and then classify each home insight embedding if the closest neighbor ($k=1$) exceeds a cosine similarity threshold α . We use α equal to 0.80 after optimizing for the precision metric defined in Table 2.

As seen in Table 1, the embedding clustering approach yields a significant increase in the total number nodes and edges in the taxonomy. However, there will still be many noisy parent to child relationships due to the disadvantage of using only textual similarity for link construction. For example, a keyphrase “golf” may become a child under the category “golf course” because the cosine similarity between these two phrase embeddings is above the α threshold, but this would not be a useful link for categorization purposes. Thus, we use the taxonomy at this stage to bootstrap training data for the link pruning algorithm described in the following section, which will aim to mitigate some of these improperly defined edges in the taxonomy.

2.2 Link Pruning and Revision

In order to remove noisy edges in the taxonomy, we set up a binary classification task to predict the probability of a valid parent-child relationship. Formally, we define a node i , its parent as $p(i)$, and an edge connecting the two nodes as $e(i, p(i))$. In addition, a candidate parent is defined as a node in the taxonomy with at least one child node.

We first manually review all edges from the taxonomy generated from Section 2.1 to remove any clearly incorrect parent-child relationships. We then use the remaining edges as positive samples for the binary classification task, yielding roughly 4000 pairs represented as $(token_i[SEP]token_{p(i)}, 1)$. The inputs $token_i$ and $token_{p(i)}$ are the tokenized

Table 2. WordNet Precision

Edge Set	Precision
Random	0.0
Embedding Similarity	0.114
Taxonomy	0.211

node tags separated by the special token $[SEP]$. 1 corresponds to the positive output label. We then create a similar number of negative samples by creating tuples $(token_i[SEP]token_{p(j)}, 0)$ where $p(j)$ is a candidate parent that does not lie on the path from a node i to the root node of the taxonomy.

Following data generation, we train a BERT-based binary text classification model on these tuples as input and remove all child nodes i from edges $e(i, p(i))$ which are classified as invalid [3]. For each pruned child node k , we create pairs with all possible candidate parents in the taxonomy and then apply the aforementioned classification model to get a probability representing the validity of each possible relationship. We then create an edge between the node k and the candidate with the maximum probability of being a parent. The last row of Table 1 displays the final state of the taxonomy after this approach.

We can use this inference protocol to automatically suggest additional edges in the taxonomy if there are new home insights on listings, which can then be manually approved or rejected by a reviewer. This human-in-the-loop step is critical for maintaining highly precise edges for the candidate recommendations described in Section 3.2 and a common feature of other taxonomy-based systems in literature [2].

3 METRICS

3.1 Quality of Taxonomy Relationships

In order to validate the quality of edges in the taxonomy, we look at all nodes in the structure that also appear in WordNet, which is a publicly available ontology linking words into semantic relations with different constructs (e.g synonyms). For each node in this subset, we evaluate three distinct approaches for edge construction (i) creating an edge between the node and a random candidate parent from the taxonomy (ii) creating an edge between the node and the candidate parent with closest embedding similarity and (iii) using the existing relationship defined in the taxonomy.

We then calculate precision by computing the number of edges from each of the approaches (i-iii) that exist in WordNet over the total number of possible WordNet edges that include a node from the taxonomy. As seen in Table 2, there is a clear improvement from building a hierarchical structure to capture more complicated WordNet relationships like hyponyms and meronyms.

3.2 Evaluation of taxonomy-based Candidate Selection for Recommendations

Table 3 compares a baseline substring match algorithm versus the taxonomy for generating candidate recommendations for top search keywords in the Seattle metropolitan area. We focus on this region due to the abundance of listings (20,000+) that are available for candidate selection.

On each property listing, “home insights” are keyphrases under descriptions that highlight interesting attribute information about the property (e.g “swimming pool”, “wood flooring”). The substring match algorithm counts a listing as a potential candidate if any substring of home insights associated with the listing matches with the search keyword

Table 3. Number of Candidate Recommendations for Top Search Keywords in Seattle

Popular Search Keyword	Baseline (substring match)	taxonomy (Parent)	taxonomy (Grandparent)
gym	162	2222 (Gyms and Studio)	3811 (Sports and Recreation)
mid century	31	98 (Vintage)	8702 (Style)
houseboat	2	103 (Boat Dock)	5806 (Waterfront)
balcony	536	733 (Loft)	1777 (Rooms)

in question. Column 2 of Table 3 displays the total number of candidate recommendations from using this approach. While this method finds a significant number of candidates for simple queries (e.g gym, balcony), it is unsurprisingly not effective for generating candidates for more complicated query terms (e.g mid century, houseboat).

Columns 3 and 4 display the total number of candidate listings from using two resolutions of categorization with the taxonomy, with the following algorithm (i) map each home insight tag on a listing to the closest category in the taxonomy, (ii) map the search keyphrase to the closest category in the taxonomy, (iii) count a listing as a candidate if there is any intersection of categories between the output from 1 and 2.

Unsurprisingly, we find a significant increase in the total number of candidate listings over the baseline approach for all top search terms in Table 3. The more interesting takeaway from this comparison is the ability to produce candidates at different resolutions of similarity with an taxonomy-based structure. For example, “gym” could be mapped to a cluster of candidate listings with home insights categorized under “Gyms and Studios” or could be mapped to a more general cluster of listings categorized under “Sports and Recreation.”

The ability to generate candidates that may only be loosely related to a search query (in terms of substring matches or direct embedding similarity) but still fall under a similar topic of interest is extremely valuable for diversifying recommendations. This is particularly crucial for a real-estate application where customers are typically unsure about what supply is available to them. Furthermore, unlike more traditional topic modeling approaches like Latent Dirichlet allocation, an taxonomy-based approach also provides a hierarchical ordering of topics that can also be dynamically adjusted based on the requirements of the recommender system.

4 CONCLUSION AND FUTURE WORK

In this extended abstract, we present an end-to-end approach for constructing a taxonomy for hierarchical recommendations of listing properties. We discuss how a hierarchical structure can capture more sophisticated edge relationships (e.g hyponyms) compared to simple embedding similarity as well as how a taxonomy can increase the total number of candidate recommendations for a real-estate search application. There are several areas of future work including (i) exploring active learning techniques to propose new edges in the taxonomy based on recently added keywords to listings, and (ii) porting the taxonomy into an ontology-based structure that could enable recommendations across child-child, child-entity, and other relationships not explored in this work.

REFERENCES

- [1] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2016. Enriching Word Vectors with Subword Information. <https://doi.org/10.48550/ARXIV.1607.04606>
- [2] Song Cui and Dhananjay Shroutry. 2020. Interest Taxonomy: A knowledge graph management system for content understanding at Pinterest. (2020). <https://medium.com/pinterest-engineering/interest-taxonomy-a-knowledge-graph-management-system-for-content-understanding-at-pinterest-a6ae75c203fd>
- [3] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. <https://doi.org/10.48550/ARXIV.1810.04805>
- [4] Bhargav Kanagal, Amr Ahmed, Sandeep Pandey, Vanja Josifovski, Jeff Yuan, and Lluís Garcia-Pueyo. 2012. Supercharging Recommender Systems using Taxonomies for Learning User Purchase Behavior. (2012). <https://doi.org/10.48550/ARXIV.1207.0136>
- [5] Jyoti Prakash Maheswari. 2019. My Internship at Zillow Group AI Part 1: Attribute Recognition in Real Estate Listings. (2019). <https://medium.com/zillow-tech-hub/my-internship-at-zillow-group-ai-part-1-attribute-recognition-in-real-estate-listings-7d2f92009552>

SPEAKER BIOS

Anish Khazane is an Applied Scientist on the Home Understanding AI team at Zillow Group. He primarily focuses on building scalable machine learning models that enable rich content understanding and home recommendations for millions of users. Prior to Zillow, he worked at Capital One on building deep neural networks for representing financial transactions for merchant recommendations and advanced language models for the company’s Eno chatbot. He holds a M.S in Computer Science from the Georgia Institute of Technology and a B.S in Computer Science from the University of California, Berkeley.

Zachary Harrison is an Applied Scientist at Zillow Group where he works on the Home Understanding AI team. His work focuses on developing machine learning models for content extraction and recommendations for Zillow’s customers. He holds a M.S. in Computer Science from the University of Massachusetts Amherst and a B.S in Computer Science and Computer Engineering from the University of Wisconsin-Madison.

A APPENDIX: ADDITIONAL EVALUATION OF TAXONOMY QUALITY

A.1 Evaluation of Subtree Similarity

Table 4. Embedding Based Subtree Similarity (add more detailed caption)

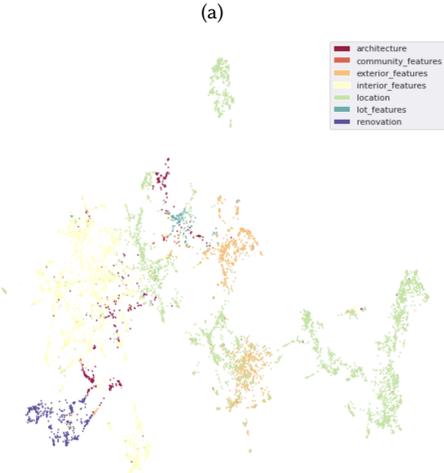
High Level Tree	High Level Tree Score	Subtree	Subtree Score	Subtree Size
Random	0.545	-	-	62
Exterior	0.658	Swimming Pool	0.778	576
Interior	0.633	Natural Light	0.755	338
Location	0.554	Dining and Drinking	0.748	214

We can plot the embedding representations of every tree node in the taxonomy to order to look for subtree clustering. Figure 1 shows this clustering of node embeddings with cosine similarity as the distance measure. This visualization technique helps with pruning the taxonomy in case there are needed hierarchy modifications, such as merging similar trees or breaking up larger ones. The location subtree, for example, is quite a large tree covering many different keywords. Furthermore, Figure 2 displays an even clearer clustering of location topics from the next level of the location subtree. We can now define a subtree similarity metric with the following protocol (i) take all unique node pairs within each subtree (ii) calculate cosine similarity between their respective embedding representations and (iii) average these similarities across all pairs in a subtree. The average similarity represents the connectivity of a subtree, with larger values indicating a more topic specific subtree.

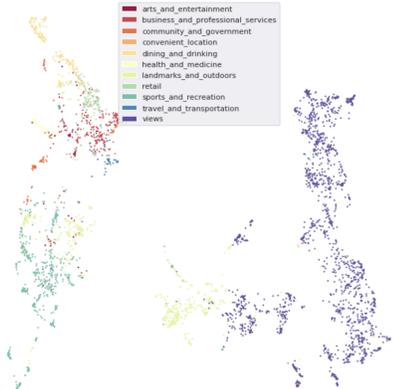
In Table 4, we display the results from computing this metric over different depths of subtrees (e.g exterior to swimming pool), with a comparison against a baseline approach of randomly creating edges between nodes and parents in the taxonomy. As seen in the table, the average similarity at all depths of different subtrees beats the random approach. We also observe that the further we move down a subtree, the higher the similarity score is which means more specific or semantically similar neighboring nodes.

We also note that the Location subtree has quite a low score, barely beating the randomly generated tree. This is due to the large size of this tree (4081) as well as the diverse topics this tree covers. For example, attractions in a city such as restaurants, museums, or sports stadiums will be quite different from attractions in a more rural setting which may include mountains, hiking trails or lakes.

A.2 Visual Representation of Taxonomy-based Recommendations



(b) Visualization of tag embeddings labeled based on their highest level subtree.



(c) Visualization of Location tag embeddings labeled based on subtrees within Location.

14 Seattle homes with golf (with taxonomy parent)

										
\$749,000	\$740,000	\$800,000	\$2,200,000	\$485,000	\$999,999	\$890,000	\$420,000	\$1,037,500	\$870,000	\$870,000
Condo, 2 Beds, 1,121 Sq.Ft., Tag: golf_course	Single Family, 4 Beds, 1,610 Sq.Ft., Tag: golf_course	Townhouse, 3 Beds, 1,491 Sq.Ft., Tag: golf_course	Single Family, 3 Beds, 2,880 Sq.Ft., Tag: golf_course	Condo, 3 Beds, 1,707 Sq.Ft., Tag: golf_course	Townhouse, 3 Beds, 1,827 Sq.Ft., Tag: golf_course	Townhouse, 3 Beds, 1,640 Sq.Ft., Tag: golf_course	Single Family, 2 Beds, 570 Sq.Ft., Tag: west_seattle_golf_course	Single Family, 4 Beds, 2,420 Sq.Ft., Tag: west_seattle_golf_course	Townhouse, 3 Beds, 1,608 Sq.Ft., Tag: golf_course	Townhouse, 3 Beds, 1,608 Sq.Ft., Tag: golf_course
Listed 23 days ago	Listed 31 days ago	Listed 23 days ago	Listed 3 days ago	Listed 3 days ago	Listed 36 days ago	Listed 36 days ago	Listed 18 days ago	Listed 17 days ago	Listed 27 days ago	Listed 27 days ago

53 Seattle homes with sports_and_recreation (with taxonomy grandparent)

										
\$829,900	\$2,200,000	\$634,000	\$1,037,500	\$624,000	\$2,700,000	\$599,000	\$904,000	\$378,000	\$378,000	\$749,000
Townhouse, 2 Beds, 1,025 Sq.Ft., Tag: volunteer_park	Single Family, 3 Beds, 2,880 Sq.Ft., Tag: golf_course	Condo, 1 Beds, 768 Sq.Ft., Tag: yoga_studio	Single Family, 4 Beds, 2,420 Sq.Ft., Tag: west_seattle_golf_course	Condo, 2 Beds, 1,285 Sq.Ft., Tag: healthy_reserves	Condo, 2 Beds, 1,890 Sq.Ft., Tag: yoga_studio	Condo, 1 Beds, 676 Sq.Ft., Tag: sky_club	Condo, 1 Beds, 751 Sq.Ft., Tag: yoga_studio	Condo, 1 Beds, 545 Sq.Ft., Tag: fitness_studio	Condo, 1 Beds, 545 Sq.Ft., Tag: fitness_studio	Condo, 2 Beds, 1,121 Sq.Ft., Tag: golf_course
Listed 23 days ago	Listed 3 days ago	Listed 151 days ago	Listed 17 days ago	Listed 8 days ago	Listed 304 days ago	Listed 65 days ago	Listed 53 days ago	Listed 2 days ago	Listed 2 days ago	Listed 23 days ago

Fig. 2. Visual example of recommendations at the parent and grandparent level of taxonomy for the search query, "golf". As seen in the image, categorizing the query under different levels of resolutions allows for a larger and more diverse number of recommendations.