



34 curation, the BioASQ initiative created a yearly challenge [1]. In this challenge, BioASQ  
35 provides millions of MeSH-labeled PubMed abstracts for competitors to train models on.  
36 Competitors can then use their models to predict MeSH labels for unlabeled abstracts, which  
37 BioASQ periodically releases in batches before they are manually indexed by the NLM.

38 Since MeSH labels are organized hierarchically, the challenge task is a variant of multilabel  
39 document classification known as hierarchical classification. A PubMed article may be  
40 indexed not just by the leaves of the MeSH DAG, but also most internal nodes. This  
41 hierarchical wrinkle is both a blessing and a curse: while it collapses the some 28,000  
42 possible labels into a more manageable space, it complicates labeling strategies. Traditional  
43 metrics for flat multilabel classification, like accuracy, are binary “hit-miss” measures which  
44 fail to capture that closely related nodes share a greater similarity than distantly related  
45 nodes.

46

47 My approach to this task can generally be outlined as follows:

48 1) Generate distributed representations of abstracts (AbsVecs).

49 2) For each MeSH term, merge related AbsVecs to generate a composite representation  
50 (TermVecs).

51 3) For unlabeled examples, generate test AbsVecs (i.e. query expansion); for each test  
52 AbsVec, rank TermVecs by their distance via some metric (e.g. cosine similarity), and  
53 threshold/thin the ranked list to remove redundant nodes (i.e. nodes on the same path). Such  
54 an approach should minimize a cost that takes the unique challenges of hierarchical  
55 classification into account.

56

## 57 **2 Background**

58

### 59 **2.1 Distributed representations of words and sentences**

60 Word-to-vec (WV) is a neural-network like architecture with a single hidden layer created by  
61 Mikolov and coworkers to generate distributed representations of words [2]. It has two  
62 flavors. The skip-gram variant (WV-SG) takes a target word as input and attempts to predict  
63 its surrounding window of words (i.e. context). The continuous-bag-of-words (CBOW)  
64 model does the opposite: it tries to predict the target word from the target word’s context  
65 (figure 2). Both models have two weight matrices (one to connect the input layer to the  
66 hidden layer and one to connect the hidden layer to the output layer). Assuming an  
67 orientation in which the rows of a weight matrix correspond to words in the vocabulary and  
68 columns correspond to elements of the hidden layer, each row of either weight matrix is a  
69 vector corresponding word in the vocabulary which, through several rounds of forward- and  
70 back-propagation, comes to capture the semantic meaning of the word. The resulting word  
71 vectors (WordVecs) can be averaged to crudely represent longer strings of text, and are also  
72 useful inputs for a variety of deep learning networks like RNNs, RNTNs, and LSTMs.

73 More recently, Mikolov and coworkers built upon WV to generate more precise vector  
74 representations of longer strings of text; this model is commonly known as paragraph-to-vec  
75 (PV) [3]. Like WV, PV comes in two variants: (PV-DM), which is analogous in structure to  
76 WV-CBOW, and PV-CBOW. These approaches are similar to their WV counterparts, except  
77 that a paragraph matrix is included. Assuming a similar orientation, rows of this matrix  
78 correspond to vector representations of documents. As the model is trained through rounds  
79 of forward- and back-propagation, the paragraph vectors come to represent how the broader  
80 document context modifies the meaning of the word vectors that comprise it. In essence,  
81 these paragraph vectors act like a document-wide “memory”.

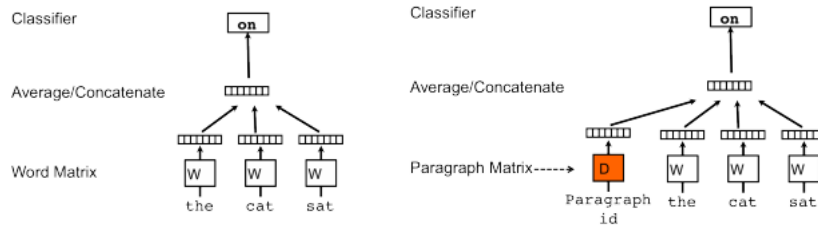


Figure 2: WV-CBOW model; PV-DM model. Images taken from Mikolov, et al. [3]

82  
83  
84  
85  
86  
87  
88  
89  
90  
91  
92  
93

## 2.2 Hierarchical clustering scoring

As mentioned in the introduction, there are unique challenges to hierarchical clustering relative to traditional clustering. For example, given a correct label of “cancer” and a predicted label of “lung-cancer”, a binary “hit-miss” metric like accuracy would give no credit to the prediction. However, a more appropriate metric would take into account that “lung-cancer” is closely related to “cancer” and view the prediction more favorably. These challenges are further discussed in detail by [4], and summarized in figure 3. Two promising metrics are MGIA and FLCA.

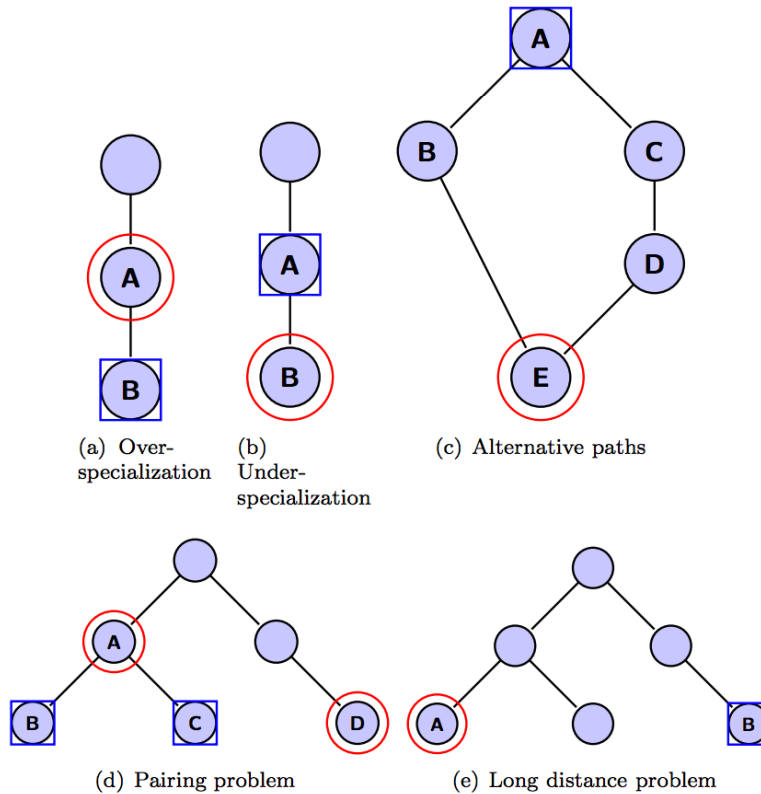


Figure 3: Five types of unique scoring challenges in hierarchical classification. While a flat classification metric would score all of these misses equally, more advanced metrics such as MGIA and FLCA would score, for example, (a) and (b) as better classifications than (e).  
Image from Kosmopoulos, et al. [4]

94  
95  
96  
97  
98  
99

## 2.3 Prior work

100 This is the fourth iteration of the BioASQ MeSH indexing task. There have been several  
 101 successful entries [1], many of which outperform the NLM's version of an automated  
 102 indexer, the Medical Terminology Indexer (MTI). Performance results on the most recent  
 103 test batch (batch 3, week 3) are shown in figure 4. The MeSHLabeler submissions, which  
 104 take advantage of article metadata, are frequently top performers [5]. See [http://participants-](http://participants-area.bioasq.org/results/4a/)  
 105 [area.bioasq.org/results/4a/](http://participants-area.bioasq.org/results/4a/) for more.

System	Acc.	LCA-F
auth1	0.4184	0.4694
auth2	0.4289	0.4690
BioASQ Filtering	0.4075	0.4634
CSX-1	0.2702	0.2964
d33p	0.1270	0.1907
Default MTI	0.4201	0.4707
iria-1	0.2567	0.3557
LABDA baseline	0.3249	0.3959
LABDA ElasticSearch	0.1819	0.3158
LargeElasticLABDA	0.1812	0.3150
MeSHLabeler	0.4747	0.5113
MeSHLabeler-2	0.4735	0.5024
MeSHLabeler-3	0.4621	0.4989
MTI First Line Index	0.4088	0.4636
UCSDLogReg	0.3420	0.3781

106 Figure 4: BioASQ Task 4a, Test batch 3, week 3 results. LCA-F is a hierarchical measure of  
 107 performance. Image from the BioASQ participants area webpage.

108  
 109 **3 Approach and experiment**

110  
 111 **3.1 Inputs**

112  
 113 **3.1.1 MeSH hierarchy**

114 The 2016 version of MeSH has some 27.9 K terms organized into a DAG. BioASQ provides  
 115 this DAG as a parent child list. Inspection of this graph through the python networkx library  
 116 revealed the properties shown in figure 5. The longest path was determined to be 17.

117

```

Type: DiGraph          Summary of shortest paths to 18714 leaves:
Number of nodes: 27900  AVG: 5.75
Number of edges: 37796  MIN: 2
Average in degree: 1.3547  MAX: 13
Average out degree: 1.3547

```

118 Figure 5: Summary of properties of the MeSH DAG.

119

120

121

**3.1.2 Dataset**

122

BioASQ provides two versions of the 2016 training set, both in JSON. The first is a 20.9 GB set with 12.2 M mesh-labeled abstracts (along with other information such as title, journal, and pubmed id). The second is a 8.5 GB subset with 4.9 M abstracts from a limited list of journals. The unlabeled test sets provided on BioASQ use the same limited list of journals. While initial efforts went into working with the complete dataset, the difficulties involved with working with the larger dataset led to the use of the limited one.

124

125

126

127

128

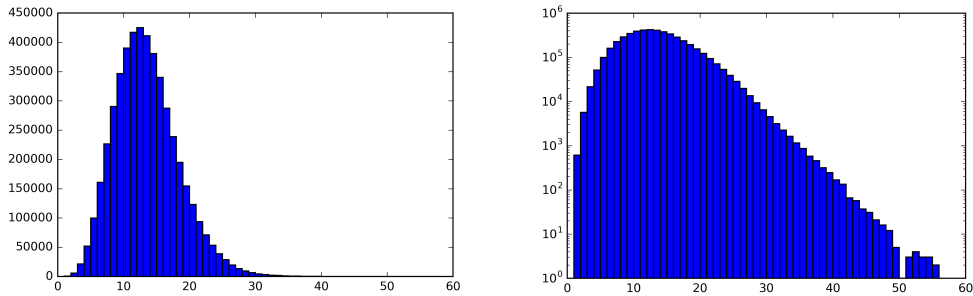
Analysis of the number of MeSH terms used to label each training example revealed an average of 13 labels with considerable variation (figures 6 and 7). This variation confounds efforts to predict labels, as the number of labels is not known in advance. Moreover, brute-force approaches to minimize cost over all possible numbers of labels

129

130

131

132



133

Figure 6: Distribution of MeSH label counts for the limited 8.5 GB dataset of abstracts (standard scale; log scale).

134

135

```

count: 4917245
min max: 1 60
avg: 13.01
med: 13.0
var std: 23.27 4.82

```

136

Figure 7: Corresponding summary statistics for figure 6 (MeSH label counts).

137

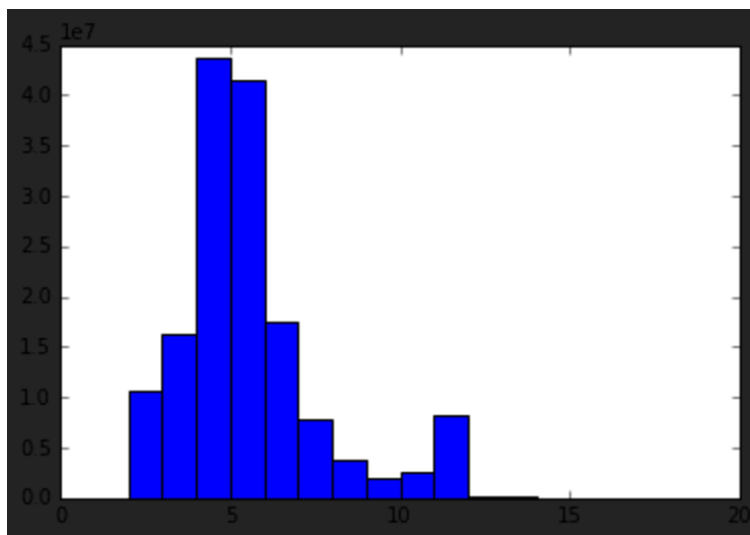
138

Additional analysis of some 154 million labels in the complete 20.9 GB dataset found a bimodal distribution of shortest paths (figures 8 and 9). This finding suggests that an intelligent search algorithm should be able to achieve higher performance relative to brute-force sampling of all 27.9 K possible MeSH terms—many of the terms are only

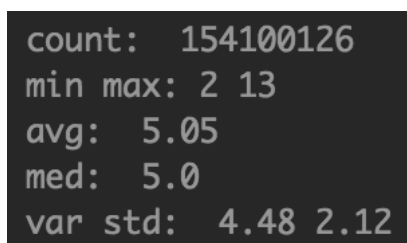
139

140

141



142 Figure 8: Distribution of shortest path lengths for all 154,100,126 labels encountered in the  
 143 complete 20.9 GB training set. The distribution appears to be bimodal.  
 144



145 Figure 9: Summary statistics for the shortest path lengths for all 154,100,126 labels  
 146 encountered in the complete 20.9 GB training set.

147  
 148  
 149

### 3.1.3 Word vectors

150 Initial word vectors (WordVecs) were provided as a resource from BioASQ and described at  
 151 <http://participants-area.bioasq.org/info/BioASQword2vec/>. Briefly, ~10.9 million abstracts  
 152 were stripped of punctuation and converted to lower case. The resulting corpus had ~9.4  
 153 million unique words. The Gensim implementation of WV was then applied to the corpus,  
 154 which yielded 1.7 million 200-dimensional WordVecs (each corresponding word had a  
 155 minimum of 5 mentions in the corpus). The corresponding 3.5 GB space-delimited file had  
 156 one vector per line; a corresponding file with a sorted list of words in the vocab indexed  
 157 these vectors.

158  
 159  
 160  
 161

## 3.2 Processing

### 3.1.1 Creating AbsVecs

162 There are many possible avenues to generating vector representations of abstracts  
 163 (AbsVecs). The simplest of these is to take a uniform average of the WordVecs for each word  
 164 in the abstract, which is ultimately the approach that was pursued in this paper. This  
 165 approach gives very coarse vectors because words that are simply mentioned more often will  
 166 tend to dominate the average, resulting in low signal-to-noise. An improvement on this  
 167 approach is to use a similar but weighted average of WordVecs. Because the task seeks to  
 168 index documents with broad terms that capture what the document is about, a term-

169 frequency inverse document-frequency (TF-IDF) weighting that effectively identifies  
170 keywords makes appropriate sense. Ultimately, an approach such as PV or RNNs to better  
171 estimate the weights on the WordVecs would likely yield even better composite  
172 representations.

173 To generate AbsVecs for each training example, raw abstract text was stripped of  
174 punctuation and converted to lowercase. The resulting words were looked up in the BioASQ  
175 provided vocab to obtain a word index, which in turn was used to locate the corresponding  
176 WordVec. The WordVecs for each word in a given abstract were then averaged to create an  
177 AbsVec (one per training example). Like the comprising WordVecs, each AbsVec was 200-  
178 dimensional.

179

### 180 **3.1.2 Creating TermVecs**

181 As with generating AbsVecs, there are many possible ways to create vector representations  
182 of MeSH terms (TermVecs). Again, the simplest possibility to generate a TermVec (and the  
183 one that was pursued in this paper) is to take all of the AbsVecs associated with a given term  
184 and to average them. This approach has a disadvantage in that it will not produce a TermVec  
185 for any term not used in the training corpus. A better approach would be to merge the  
186 average of a term's AbsVecs the average of its children's TermVecs. A simple 50-50  
187 weighting structure, would be the first thing to try here although a neural network could also.  
188 Indeed, the hierarchical structure of the MeSH resembles that of a Recursive Neural  
189 Network, so it appears likely that an RNN-like arch.

190 Like the comprising AbsVecs, each TermVec was 200-dimensional.

191

### 192 **3.1.3 Predicting MeSH labels**

193 Regrettably, while I was able to produce a basic set of training AbsVecs and TermVecs, I  
194 failed to generate a set of validation or test AbsVecs within the constraints of the project  
195 timeline. Given additional time, I would have liked to have *a*) performed 3-fold cross-  
196 validation when generating TermVecs from the training AbsVecs (such code to generate 3-  
197 fold TermVecs from AbsVecs is included in the submission), and *b*) generated AbsVecs from  
198 unlabeled test data to make predictions for submission to BioASQ.

199 To make a prediction for an unlabeled test AbsVec, I would have calculated a list of  
200 TermVecs ranked and sorted on their distance to the AbsVec (e.g. via cosine similarity or  
201 Euclidean distance). If necessary, a threshold could be applied to the ranked list (e.g. top  
202 10%). Then, given a fixed number of labels, I would have attempted to find the pruning of  
203 that ranked list that gave the highest combined while enforcing the NLM-dictated constraint  
204 that any predicted term have no predecessors (ancestors) that are also predictions (e.g.  
205 "cancer" would not be a valid label if "lung-cancer" was a superior prediction). In the case  
206 of validation, the cost of the predicted terms could be established using accuracy (flat  
207 metric) and F-LCA (hierarchical metric).

208

## 209 **4 Conclusions**

210 Progress on the project was greatly hampered by the size of the datasets involved (20.9 GB  
211 for the complete dataset, 8.5 GB for the 3.5 GB for the WordVectors). Downloading and  
212 decompressing the data proved alone proved to be nontrivial, as BioASQ used a compression  
213 protocol not documented on their website. Over the course of the project, a significant  
214 amount of time was also invested into working with the large 20.9 GB dataset. After  
215 struggling with the size of the dataset (which would not fit in RAM), an effort was made to  
216 organize data and resulting calculations into a sqlite database for faster querying; however,  
217 while the database was successfully created, this amounted to a dead end effort with little  
218 useful speedup, and ultimately the decision was made to use the smaller, more limited 8.5  
219 GB dataset. Future work would greatly benefit from better organization and access to the  
220 large datasets.

221 Nonetheless, the work above has produced a set of AbsVecs (each the average of comprising

222 WordVecs) and TermVecs (each the average of comprising AbsVecs) for the limited training  
223 set data provided by BioASQ. The usefulness of these two resources remains to be  
224 evaluated. As mentioned, both the generation of the AbsVecs and the TermVecs could easily  
225 be improved beyond the by using smarter weighting (e.g. by using TF-IDF for AbsVecs or  
226 propagation of TermVecs through the DAG) or by using a neural architecture to learn  
227 appropriate weights. More generally, an integrated learning model that was trained to give  
228 task specific (rather than general purpose) AbsVecs and TermVecs based directly on  
229 predictive performance (e.g. measured by F-LCA) is also likely to give better results;  
230 however, it remains unclear to me how to adapt deep networks to the unique challenges of  
231 multilabel hierarchical classification (at least in an efficient way that takes advantage of the  
232 hierarchical structure).

233 There are several additional possibilities for improvement. The MeSHLearner model  
234 established that other abstract metadata can providing invaluable features for better  
235 classification [5]. This finding makes tremendous sense—for instance, certain journals (e.g.  
236 “Cancer”) are likely to have very common MeSH annotations. Rather than comparing  
237 representations of abstracts against representations of terms, it may ultimately be more  
238 helpful to conceptualize embedding abstracts as some representation of a DAG and attempt  
239 to reconcile the predicted DAG with some pruning of the full MeSH DAG structure  
240 (inspiration from [6]).

#### 241 **Acknowledgments**

242 The author thanks the CS224d teaching assistants for their hard work this quarter as well as  
243 Dr. Socher for his excellent lectures and guidance.

#### 244 **References**

- 245 [1] G. Tsatsaronis, et al: *An overview of the BIOASQ large-scale biomedical semantic indexing and*  
246 *question answering competition*. Apr. 30, 2015. BMC Bioinformatics.
- 247 [2] T. Mikolov, et al: *Efficient estimation of word representations in vector space*. Version 3. Sept. 7,  
248 2013. arXiv:1301.3781v3. [cs.CL]
- 249 [3] Q. Le & T. Mikolov: *Distributed representations of sentences and documents*. Version 2. May 22,  
250 2014. arXiv:1405.4053v2 [cs.CL]
- 251 [4] A. Kosmopoulos, et al: *Evaluation measures for hierarchical classification: a unified view and*  
252 *novel approaches*. Version 2. Jul. 1, 2013. arXiv: 1306.6802v2 [cs.AI]
- 253 [5] K. Liu, et al: *MeSHLabeler: improving the accuracy of large-scale MeSH indexing by integrating*  
254 *diverse evidence*. Jun., 2015. Bioinformatics.
- 255 [6] S. Bengio, J. Weston, & D. Grangier: *Label embedding trees for large multi-class tasks*. 2010.  
256 *Advances in Neural Information Processing Systems* 27, 163-171.