# Understanding the Stability of Medical Concept Embeddings

Grace E. Lee
School of Computer Science and Engineering
Nanyang Technological University, Singapore
leee0020@e.ntu.edu.sg

Aixin Sun
School of Computer Science and Engineering
Nanyang Technological University, Singapore
axsun@ntu.edu.sg

## ABSTRACT

Frequency is one of the major factors for training quality word embeddings. Several work has recently discussed the stability of word embeddings in general domain and suggested factors influencing the stability. In this work, we conduct a detailed analysis on the stability of concept embeddings in medical domain, particularly the relation with concept frequency. The analysis reveals the surprising high stability of low-frequency concepts: low-frequency (<100) concepts have the same high stability as high-frequency (>1000) concepts. To develop a deeper understanding of this finding, we propose a new factor, *the noisiness of context words*, which influences the stability of medical concept embeddings, regardless of frequency. We evaluate the proposed factor by showing the linear correlation with the stability of medical concept embeddings. The correlations are clear and consistent with various groups of medical concepts. Based on the linear relations, we make suggestions on ways to adjust the noisiness of context words for the improvement of stability. Finally, we demonstrate that the proposed factor extends to the word embedding stability in general domain.

## KEYWORDS

Medical concept embedding, Stability, Context words, Frequency, UMLS

## 1 INTRODUCTION

Medical concepts are medical terminologies linked to Unified Medical Language System (UMLS)[1]. They are frequently used in medical domain for accurate and effective communication. Using medical concepts identified by unique IDs (*e.g.,* CUI in UMLS) helps avoid misunderstanding and allows to utilize semantic information in UMLS.

Word embeddings encode syntactic and semantic aspects of words into low-dimensional and dense representations. As they have shown great advance, many studies have proposed approaches to learning medical concept embeddings using various biomedical resources.

Several studies have shown the instability of word embeddings in general domain [2, 13, 14, 25]. Whenever word embeddings are

---

[1]UMLS is the metathesaurus in medical domain which provides an unified language system over various medical dictionaries such as ICD-10, MeSH, and SNOMED-CT. More details are in https://www.nlm.nih.gov/research/umls/

trained on the same training corpus with the same parameter setting, encoded relations among words should be consistent in each time. It is reported however that word embeddings fail to learn consistent relations, shown by using a stability measure (more details later).

The instability of embeddings poses greater problems in medical domain. Unstable medical concept embeddings result in a wide-range of negative consequences from unsuccessful document filtering to misdiagnosis. While medical concept embeddings are of importance, their stability has not been carefully investigated. In this work, we focus on medical concept embeddings and examine their stability. We use 'word embedding' and '(medical) concept embedding' interchangeably throughout this work.

It is a well-known premise that frequency is one of the major factors contributing to quality word embeddings. As word embeddings are trained using word co-occurrence information, a high word frequency in a corpus provides sufficient training instances and thus helps obtain a quality word embedding. We conduct a detailed analysis on the relation between the stability of medical concept embeddings and concept frequency. The analysis reveals that some low-frequency concepts result in high stability as high-frequency concepts do. For example, the frequencies of two concepts, Anaesthesia procedure and Knee joint operation are 10,377 and 32 respectively in our dataset and both concepts have the same high level of stability 0.8. Our analysis results show that there are other factors that influence the stability of medical concept embeddings in addition to frequency.

Finally, we propose a new factor, *the noisiness of context words*, which affects the stability of medical concept embeddings. The rationale is: if the context words of concept are noisy, its embedding undergoes inconsistent training in various directions, leading to the low quality of embedding. We use a measure called *normalized entropy* to compute the noisiness of context words. It quantifies how flat or peaked a distribution of elements (context words in this context) is. When a word has a high frequency in a corpus, its number of context words is also high. Normalized entropy negates the effect of the total number of context words using normalization. In this way, we measure solely the noisiness of context words without taking frequency into account.

In medical domain word2vec is widely adopted when training medical concepts embeddings. We focus on word2vec [18] as a representative among existing word embedding algorithms. In the experiments, we demonstrate the relation between the proposed factor and the stability of concept embeddings by calculating a linear correlation coefficient. The evaluation result shows moderate correlations. The correlations are clear and consistent with various groups of medical concepts. The proposed factor provides an empirical reasoning of how low-frequency concepts are able to result in high-stability embeddings, despite the small number

of training instances. Besides, based on the moderate linear cor-relation, we make suggestions on ways to adjust the noisiness of context words. Lastly, we extend the evaluation to general-domain word embeddings and show the equivalent linear relation of the proposed factor.

## 2 PRELIMINARIES AND RELATED WORK

In this section, we first review word embedding algorithms, and then introduce training medical concept embeddings with different types of medical data. Next, we provide an overview of previous studies on the stability of word embeddings.

### 2.1 Word Embeddings

Word embeddings are a fundamental building block in NLP tasks. It represents syntactic and semantic meanings of a word into a low-dimensional vector of real values. There are several word em-bedding algorithms such as word2vec [18], GloVe [21], and Fast-Text [5]. While details in training word embeddings are different in algorithms, they share the same hypothesis: words that occur in the same contexts tend to have similar meanings.

In this paper, we focus on word2vec, particularly the skip-gram model with negative sampling [18], as it is a popular model to train embeddings in medical domain. It trains an embedding by using co-occurrences between a given word and its neighboring words (context words) that appear within a pre-defined window size. The objective of skip-gram model is to predict the context words given target word. Initially, word embeddings are randomly initialized. During the training, the skip-gram model find a set of word embeddings, which maximizes the objective function:

$$\frac{1}{T} \sum_{i=1}^{T} \sum_{-k \leq j \leq k, \, j \neq 0} \log P\left(w_{i+j} | w_i\right) \qquad (1)$$

where $k$ is a window size, and $w_{i+j}$ indicates a context word for a given target word $w_i$ in distance of $j$. The sign of $j$ indicates a preceding or following direction of $w_i$. Whenever word embeddings are trained, the absolute values in word embeddings are different due to the random initialization of word embeddings by design. In other words, they are in different embedding spaces. However, they are supposed to encode the consistent meanings of words because they are trained on the same training corpus with the same parameter setting.

Contextual word embeddings (*e.g.,* BERT, ELMo) have advanced word representations in recent years. They necessitate computa-tional power and large amounts of data to train and/or fine-tune thousands of parameters. In medical domain, these conditions are challenging to meet because a large scale of medical data (*e.g.,* clin-ical records) is often restricted to the public over privacy concerns. Thus, word embeddings are still frequently used tools to start off with. They are computationally inexpensive and fast, leading to faster prototyping and development. Moreover, word embeddings show effective performance even when a training corpus is small. Previous literature reports that using a small and topic-specific corpus is more effective than using a large and general corpus for domain-specific word embeddings [1, 8, 27].

Table 1: Examples of input text for training concept embed-dings. Medical concepts are presented in square brackets. A word embedding algorithm such as word2vec is applied.

| | |
|---|---|
| **Original sentence** | Calcium carbonate appears to be as effective as alu-minum hydroxide in binding dietary phosphorus in hemodialysis patients. |
| **cui2vec** | [calcium_carbonate] appears to be as effective as [aluminum_hydroxide] in binding dietary [phos-phorus] in [hemodialysis] patients. |
| **NLM** | [calcium_carbonate] [aluminum_hydroxide] [phosphorus] [hemodialysis] |

### 2.2 Medical Concept Embeddings

Medical concepts (*i.e.,* medical terminologies linked to a knowledge base UMLS) are so ubiquitous that they are in diverse medical applications/resources, such as electronic health records, health insurance claims, and biomedical literature. Each medical resource provides additional information specialized in its own application, in addition to medical concepts. Previous studies have developed various approaches for effective learning of concept embeddings by leveraging extra information given by different resources.

Electronic health records (EHRs) are patients' visit records to hospitals. Each record consists of medical concepts related to given visit, as well as patient's demographic information and medical history. Records are also tagged with timestamps of visits. EHRs enable effective learning of concept embeddings by incorporating information shared by simialr patients or temporal information, in addition to co-occurrences of concepts within records [6, 9, 10]. Health insurance claims are similar to EHRs and each claim includes medical concepts tagged with temporal information. However, the occurrences of claims tend to be sparse and irregular compared to EHRs. Sporadic claims make it challenging to learn the relatedness between medical concepts into embeddings. In [11], the authors have introduced grouping and shuffling techniques to mitigate the challenge.

Lastly, biomedical literature is in formal written language and there is a significant body of work to train concept embeddings on biomedical literature [17, 20, 26]. Unlike the previous two resources, it does not have pre-tagged medical concepts in it. To train concept embeddings, concepts must be identified first by using off-the-shelf tools such as MetaMap [3], QuickUMLS [23]. The most simple and straightforward approach using biomedical literature is called (**cui2vec**) [4]. It considers each concept a single word and apply a word embedding algorithm to concepts and their surrounding non-medical words. Another similar approach (**NLM**) is to apply an algorithm to only concepts, eliminating non-medical words [12]. Table 1 presents example input texts of the two approaches. A word embedding algorithm is then applied to the input text to train concept embeddings.

Among a wide range of approaches for medical concept em-beddings, in this work we use **cui2vec** and **NLM** to train medical concept embeddings. They are general enough to be adopted not only in biomedical literature but also in EHRs and insurance claims.

For example, NLM can be directly applied for EHRs and insurance claims where each record/claim consists of a bag of medical concepts. Lastly, while EHRs and health insurance claims are often privately owned, biomedical literature is publicly available.

## 2.3 Instability of Word Embeddings

In recent years, there have been several studies reporting the instability of word embeddings in general domain. It is shown that when word embeddings are trained with the same corpus and hyperparameters but with random weight initialization in multiple times, the relations among words, in particular the nearest neighbor words, are different across the sets of word embeddings. Thus, the word embeddings are considered unstable.

Existing studies have suggested several factors that influence the stability of word embeddings. They are classified into two groups: corpus-level factors and word-level factors. As corpus-level factors, in [2], the authors explore various values of a corpus size and different document lengths in the corpus. In [1, 25, 27], it is shown that a small and topic-specific corpus helps train more stable word embeddings than a large and general corpus.

As word-level factors, word frequency, part-of-speech (POS) tags, and concreteness of word's meaning are studied [22, 25]. In [25] the authors discover that words in some POS tags tend to have more stable embeddings than other tags. They also report that frequency is not a major factor on the stability of word embeddings. However, the authors do not further examine the reasons behind the finding.

Across corpus-level and word-level factors, most factors are related to word frequency either directly or indirectly. For example, in corpus-level factors, different sizes of corpus and topic specificity in a corpus incur the changes of word frequency. In word-level factors, part-of-speech tags are also related to word frequency, since some part-of-speech tags appear more frequently than other tags in a corpus.

## 3 ANALYSIS

Frequency is one of the well-known factors that influence the quality of word embeddings. Fundamentally word embedding algorithms (*e.g.,* word2vec and GloVe) train word embeddings using co-occurrence information of words. A high frequency of a word in a corpus is equivelent to a large number of co-occurred words (*i.e.,* training instances), so that it helps produce a quality word embedding. In this section, we conduct a data180led analysis on the stability of medical concept embeddings with concept frequency in a training corpus. We first introduce the stability measure and describe the settings used in training medical concept embeddings. Then, we discuss the analysis.

## 3.1 Stability Measure

We use the stability measure to evaluate the quality of medical concept embeddings [2, 25]. It is initially proposed for word embedding in general domain. We first explain the stability in the context of general-domain word embeddings. Next, we present how it is applied to medical concept embeddings.

A stability value represents how stable a word embedding is. It is the portion of overlapping words between the $n$ nearest neighbors from different embedding spaces. Each embedding space is obtained by training word embeddings with the same training corpus and hyperparameters but with random initializations of embeddings and random sampling in a negative sampling technique, if it is used in word2vec.

Formally, given a word $w$ and two embedding spaces, $P$ and $Q$, let $P_w$ and $Q_w$ be the $n$ nearest neighbors of $w$ in $P$ and $Q$ based on embedding similarities, respectively. The stability value for word $w$'s embedding is the ratio of overlapping words in $P_w$ and $Q_w$.

$$stability(w) = \frac{|P_w \cap Q_w|}{n} \qquad (2)$$

If $P_w$ and $Q_w$ consist of the same set of words, $stability(w) = 1.0$. It demonstrates the embedding of $w$ encodes consistent semantic information of $w$. Hence, $w$'s embedding is considered stable.

In our problem setting, we replace a word $w$ with a medical concept $c$. The stability of concept embeddings is defined as the portion of overlapping nearest concepts from the three spaces $P$, $Q$, and $R$, *i.e.,* $\frac{|P_c \cap Q_c \cap R_c|}{n}$, instead of two spaces. In doing so, the stability of medical concept embeddings is evaluated by stricter conditions, since often biomedical applications require a high level of accuracy. The number of the nearest neighbors is set to 10 ($n$=10), and the similarity between concept embeddings is calculated by cosine similarity [25].

## 3.2 Training Medical Concept Embeddings

There have been many studies on approaches for training medical concept embeddings with various biomedical data (details in Related Work). In this work, we examine two approaches, cui2vec [4] and NLM [12], using biomedical literature data.

**Dataset.** We use OHSUMED dataset to train medical concept embeddings. It consists of 348,566 abstracts sampled from a MEDLINE corpus. Though a MEDLINE corpus is much larger, multiple studies report that compared to a large and general corpus, a small and topic-specific corpus is more effective for embeddings in domain-specific applications [8, 27]. For medical concept identification, there are several tools available such as MetaMap [3], NCBO BioPortal [19], PubTator [24], and QuickUMLS [23]. We use QuickUMLS because it is faster than other tools. Total 40,625 (unique) medical concepts are extracted.

The vast majority of concepts have low frequencies. Figure 1 shows a comparison between the frequency distributions of words (left) and the extracted concepts (right). As expected, word frequencies follow a power law distribution. The distribution of concept frequencies also shows that the very similar distribution of words. When we consider that a high frequency helps train quality word embedding, most concepts are likely to have low-quality embeddings due to their low frequency.

Finally, the original sentences are transformed with medical concepts, as shown in Table 1. A word embedding algorithm is applied and trains concept embeddings.

**Word2vec setting.** We use the skip-gram model with negative sampling implemented in Gensim[2]. The embedding size is set to 200 and the minimum frequency is set to 5. The window size for context words and the number of epochs are set to 7 and 50, respectively. The remaining hyperparameters are set to default values. After rare

---

[2]https://radimrehurek.com/gensim/index.html

**Figure 1: Distributions of words (left) and medical concepts (right) over frequencies in the OHSUMED dataset. The majority of concepts have low frequency.**

**All concepts, $\rho$=0.189** (cui2vec)

| Stability \ Medical concept frequency bins | [5,10) | [10, 100) | [100, 1K) | [1K, 10K) | [10K, 100K) | [100K, 1M) |
|---|---|---|---|---|---|---|
| 1.0 | 2 | 52 | 196 | 237 | 52 | 0 |
| 0.9 | 47 | 515 | 1131 | 584 | 84 | 2 |
| 0.8 | 181 | 1651 | 1824 | 527 | 46 | 1 |
| 0.7 | 449 | 2652 | 1495 | 233 | 5 | 0 |
| 0.6 | 730 | 2773 | 870 | 90 | 6 | 0 |
| 0.5 | 925 | 2326 | 439 | 30 | 2 | 0 |
| 0.4 | 937 | 1569 | 184 | 14 | 0 | 0 |
| 0.3 | 605 | 895 | 62 | 2 | 0 | 0 |
| 0.2 | 367 | 376 | 22 | 0 | 0 | 0 |
| 0.1 | 131 | 126 | 3 | 0 | 0 | 0 |
| 0.0 | 22 | 19 | 0 | 0 | 0 | 0 |

**All concepts, $\rho$=0.200** (NLM)

| Stability \ Medical concept frequency bins | [5,10) | [10, 100) | [100, 1K) | [1K, 10K) | [10K, 100K) | [100K, 1M) |
|---|---|---|---|---|---|---|
| 1.0 | 4 | 29 | 113 | 163 | 46 | 1 |
| 0.9 | 42 | 419 | 768 | 527 | 87 | 2 |
| 0.8 | 243 | 1576 | 1614 | 541 | 47 | 0 |
| 0.7 | 611 | 2680 | 1636 | 289 | 7 | 0 |
| 0.6 | 962 | 3012 | 1113 | 85 | 2 | 0 |
| 0.5 | 1151 | 2371 | 526 | 18 | 1 | 0 |
| 0.4 | 921 | 1514 | 177 | 1 | 0 | 0 |
| 0.3 | 625 | 728 | 63 | 1 | 0 | 0 |
| 0.2 | 244 | 316 | 19 | 0 | 0 | 0 |
| 0.1 | 77 | 96 | 1 | 0 | 0 | 0 |
| 0.0 | 12 | 10 | 0 | 0 | 0 | 0 |

**Figure 2: Distributions of medical concepts over frequency and the stability of embeddings trained by cui2vec and NLM in heatmaps. Total number of concepts is 25,491. Each cell indicates the number of medical concepts in the corresponding bin intersected by frequency and stability.**

words are excluded using the minimum frequency threshold, total 25,491 concept embeddings are trained.

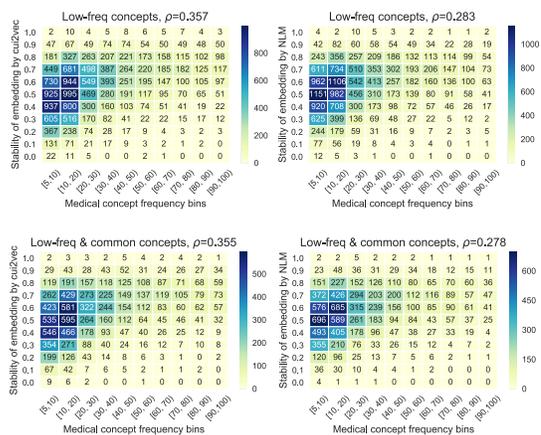Different hyperparameter values affect the quality of word embeddings [7, 15, 18]. We explore different values for five hyperparameters: the context window size, epochs, the number of negative samples, the smoothing parameter for the negative sample distribution, and the subsampling rate. We examine the changes of embedding stability when different values are used in the training. However, the evaluation results show that they result in marginal changes and the average stabilities are very similar. We report the results in Appendix.

## 3.3 Analysis: Stability of Medical Concept Embeddings and Frequency

So far, we have introduced the stability measure and the definition of stability used in this work. We have also described the setting in training concept embeddings. Now, we present a detailed analysis on the stability of concept embeddings with frequency and discuss their relations.

Figure 2 shows a distribution of medical concept embeddings over frequency and embedding stability, trained by cui2vec and NLM. In a subfigure, the $x$-axis indicates the ranges of concept frequency, and the $y$-axis indicates the stability from 0.0 to 1.0. Each cell denotes the number of medical concepts that belong to the corresponding frequency and stability bins. The color of cells shows the relative number of concepts.

As shown in Figure 2, high-frequency concepts tend to have high-stability concept embeddings. When a frequency is greater than 1,000, most concepts have stability higher than 0.5. The result confirms that indeed high frequency helps produce high stability

**Low-freq concepts, $\rho$=0.357** (cui2vec)

| Stability \ Medical concept frequency bins | [5,10) | [10,20) | [20,30) | [30,40) | [40,50) | [50,60) | [60,70) | [70,80) | [80,90) | [90,100) |
|---|---|---|---|---|---|---|---|---|---|---|
| 1.0 | 2 | 10 | 4 | 5 | 8 | 6 | 5 | 7 | 4 | 7 |
| 0.9 | 47 | 67 | 49 | 74 | 74 | 54 | 50 | 49 | 48 | 50 |
| 0.8 | 181 | 327 | 263 | 207 | 221 | 173 | 158 | 115 | 102 | 98 |
| 0.7 | 449 | 681 | 408 | 387 | 264 | 220 | 185 | 182 | 125 | 117 |
| 0.6 | 730 | 944 | 549 | 393 | 251 | 195 | 147 | 100 | 105 | 97 |
| 0.5 | 925 | 995 | 469 | 280 | 191 | 117 | 95 | 70 | 65 | 51 |
| 0.4 | 937 | 800 | 300 | 160 | 103 | 74 | 51 | 41 | 19 | 22 |
| 0.3 | 600 | 510 | 170 | 82 | 41 | 22 | 22 | 15 | 17 | 12 |
| 0.2 | 367 | 238 | 74 | 27 | 17 | 9 | 4 | 3 | 2 | 3 |
| 0.1 | 131 | 71 | 21 | 17 | 9 | 3 | 2 | 1 | 2 | 0 |
| 0.0 | 22 | 11 | 5 | 0 | 0 | 2 | 1 | 0 | 0 | 0 |

**Low-freq concepts, $\rho$=0.283** (NLM)

| Stability \ Medical concept frequency bins | [5,10) | [10,20) | [20,30) | [30,40) | [40,50) | [50,60) | [60,70) | [70,80) | [80,90) | [90,100) |
|---|---|---|---|---|---|---|---|---|---|---|
| 1.0 | 4 | 2 | 10 | 5 | 3 | 2 | 2 | 1 | 1 | 2 |
| 0.9 | 42 | 82 | 60 | 58 | 54 | 49 | 34 | 22 | 28 | 19 |
| 0.8 | 243 | 356 | 257 | 209 | 186 | 132 | 113 | 114 | 99 | 54 |
| 0.7 | 611 | 734 | 510 | 353 | 302 | 193 | 206 | 147 | 104 | 73 |
| 0.6 | 962 | 1106 | 542 | 413 | 257 | 182 | 160 | 136 | 100 | 63 |
| 0.5 | 1151 | 982 | 456 | 310 | 173 | 139 | 80 | 91 | 58 | 41 |
| 0.4 | 920 | 708 | 300 | 173 | 98 | 72 | 57 | 46 | 26 | 17 |
| 0.3 | 628 | 399 | 136 | 69 | 48 | 27 | 22 | 5 | 12 | 2 |
| 0.2 | 244 | 179 | 59 | 31 | 16 | 9 | 7 | 2 | 3 | 5 |
| 0.1 | 77 | 56 | 19 | 8 | 4 | 3 | 4 | 0 | 1 | 1 |
| 0.0 | 12 | 5 | 3 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |

**Low-freq & common concepts, $\rho$=0.355** (cui2vec)

| Stability \ Medical concept frequency bins | [5,10) | [10,20) | [20,30) | [30,40) | [40,50) | [50,60) | [60,70) | [70,80) | [80,90) | [90,100) |
|---|---|---|---|---|---|---|---|---|---|---|
| 1.0 | 2 | 3 | 4 | 2 | 4 | 2 | 4 | 2 | 1 | |
| 0.9 | 23 | 43 | 28 | 43 | 52 | 31 | 24 | 26 | 27 | 34 |
| 0.8 | 119 | 191 | 157 | 118 | 125 | 108 | 87 | 71 | 68 | 59 |
| 0.7 | 262 | 420 | 273 | 225 | 149 | 137 | 119 | 105 | 79 | 73 |
| 0.6 | 423 | 581 | 322 | 244 | 154 | 112 | 83 | 60 | 62 | 57 |
| 0.5 | 535 | 595 | 264 | 160 | 112 | 64 | 45 | 46 | 41 | 32 |
| 0.4 | 546 | 466 | 178 | 93 | 47 | 40 | 26 | 25 | 12 | 9 |
| 0.3 | 354 | 271 | 88 | 40 | 24 | 16 | 12 | 7 | 10 | 8 |
| 0.2 | 199 | 126 | 43 | 14 | 8 | 6 | 3 | 1 | 0 | 2 |
| 0.1 | 67 | 42 | 7 | 6 | 5 | 2 | 1 | 1 | 2 | 0 |
| 0.0 | 9 | 6 | 2 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |

**Low-freq & common concepts, $\rho$=0.278** (NLM)

| Stability \ Medical concept frequency bins | [5,10) | [10,20) | [20,30) | [30,40) | [40,50) | [50,60) | [60,70) | [70,80) | [80,90) | [90,100) |
|---|---|---|---|---|---|---|---|---|---|---|
| 1.0 | 2 | 1 | 2 | 1 | 1 | 1 | 1 | 1 | 1 | 1 |
| 0.9 | 23 | 48 | 36 | 31 | 29 | 34 | 18 | 12 | 15 | 11 |
| 0.8 | 151 | 227 | 152 | 126 | 110 | 80 | 65 | 70 | 60 | 36 |
| 0.7 | 372 | 426 | 294 | 203 | 200 | 112 | 116 | 89 | 57 | 47 |
| 0.6 | 576 | 685 | 315 | 239 | 156 | 100 | 85 | 90 | 61 | 41 |
| 0.5 | 696 | 589 | 261 | 183 | 94 | 84 | 43 | 57 | 37 | 25 |
| 0.4 | 493 | 405 | 178 | 96 | 47 | 38 | 27 | 33 | 19 | 4 |
| 0.3 | 355 | 210 | 76 | 33 | 26 | 15 | 12 | 4 | 7 | 2 |
| 0.2 | 120 | 96 | 25 | 13 | 7 | 5 | 6 | 2 | 1 | 1 |
| 0.1 | 36 | 30 | 10 | 4 | 4 | 1 | 2 | 0 | 0 | 0 |
| 0.0 | 4 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

**Figure 3: Distributions of low-frequency concepts (in the first row) and commonly used concepts among low-frequency concepts (in the second row) over frequency and the stability of embeddings.**

embeddings with the sufficient training instances. Overall, cui2vec and NLM show similar distributional patterns.

We observe interesting stability distributions among low-frequency concepts. The stability of low-frequency concepts (frequency < 100) varies widely from 0.0 to 1.0. They end up with the high stability embeddings even with the much smaller number of training instances than high-frequency concepts. For example, here are two medical concepts, Anaesthesia procedure (CUI: C0002903) and Knee joint operation (CUI: C0187769), with drastically different frequencies. Anaesthesia procedure has top 10% frequency (frequency: 10,377) and Knee joint operation has bottom 10% frequency (frequency: 32) in the corpus. However, both result in the very high stability 0.8 in cui2vec and NLM.

We take a close look at low-frequency concepts. We zoom into the first two columns in Figure 2, and present the distributions of low-frequency concepts alone in Figure 3. We also provide the distributions of commonly used concepts among low-frequency concepts in Figure 3. In medical domain, concepts of diseases, symptoms, treatments, and diagnostic tests are commonly looked up by medical professionals [16]. Among all concepts (Figure 2), a large portion of concepts has a frequency lower than 100, and they are also important concepts (Figure 3). In Figure 3, more than a half of concepts have the stability greater than 0.5. Some low-frequency concepts have the stability even greater than 0.8.

We compute Pearson correlation coefficient ($\rho$) between concept frequency and embedding stability of all concepts. The correlation coefficient is 0.189 and 0.200 in cui2vec and NLM, respectively. Frequency shows a weak linear correlation with the embedding stability, because low-frequency concepts also have high stability as high-frequency concepts do.

**Summary.** In this section, we examine the stability of medical concept embeddings with frequency. High-frequency concepts have high-stability embeddings, showing that a high frequency indeed helps obtain quality embeddings with sufficient training instances. More importantly, we observe that surprisingly low-frequency

concepts also have high-stability embeddings, despite the smaller number of training instances. The high-stability embeddings of low-frequency concepts are not aligned with the common understanding of the impact of frequency. This results demonstrate that there might be other factors influencing the stability of concept embeddings in addition to frequency.

## 4 PROPOSED MEASURE

We propose a new factor, *the noisiness of context words*, which influences the stability of concept embeddings. Our hypothesis is as follows: if a concept has noisy context words, its embedding will undergo inconsistent learning, leading to a low-quality embedding. Conversely, if a concept has coherent context words, its embedding will go through consistent training, and thus it results in a high-quality embedding. We note that the proposed factor focuses on the distribution of context words only without considering the total number of context words, which is dependent on frequency. In the following discussion, we use 'concept' and 'word' interchangeably. The idea is the same.

To measure the noisiness of context words, we use *normalized entropy*. Normalized entropy quantifies a distribution of elements (*i.e.,* context words in our context). It negates the effect of the total number of context words in the estimate using normalization. Formally, for a given word $w$, normalized entropy $H(w)$ is computed as follows.

$$H(w) = - \sum_{w_i \in C_w} \frac{P(w_i) \cdot \log P(w_i)}{\log |C_w|} \qquad (3)$$

In the above equation, $C_w$ denotes a set of context words that co-occur with word $w$ within the window size. $P(w_i)$ is the relative frequency of context word $w_i$ in the collection of all context words. If context words of a target word are evenly distributed (flat distribution), normalized entropy is high; if a few context words account for a large portion in the entire context words (skewed distribution), the value is low.

We evaluate the proposed factor by computing a linear correlation coefficient (Pearson correlation coefficient) with the stability of concept embeddings. Table 2 presents coefficient values calculated in three groups of medical concepts: all concepts, low-frequency (frequency < 100) concepts , and commonly used concepts among low-frequency concepts. For all concepts, the noisiness of context words shows a moderate linear correlation with the stability of concept embeddings. In general, the correlation is stronger in cui2vec than in NLM.

As shown in Figures 2 and 3, while high-frequency concepts tend to have high stability, low-frequency concepts show varied embedding stabilities. To better understand the surprising stability of low-frequency concepts, we look into low-frequency concepts and conduct a focused evaluation on them.

In Table 2, the proposed factor consistently shows a moderate linear correlation with the stability of low-frequency concepts. The correlation in cui2vec is even comparable to the correlation estimated for all concepts. In NLM, the correlation is slightly decreased, compared to the correlation for all concepts. The similar correlations are observed for commonly used concepts as well. Lastly, the proposed factor shows much stronger correlations with the stability of embeddings, compared to frequency as shown in Figures 2 and 3.

Table 2: Pearson correlation coefficient of the noisiness of context words with the stability of medical concepts. We present three groups of medical concepts: all concepts, concepts with low-frequency (<100) (Low-freq), and commonly used concepts in biomedical applications among low-frequency concepts (Low-freq & Common). All correlation coefficients are statistically significant at $p < 0.001$.

|  | All concepts | Low-freq | Low-freq & Common |
|---|---|---|---|
| cui2vec | -0.593 | -0.540 | -0.549 |
| NLM | -0.477 | -0.370 | -0.329 |

All correlation coefficients are statistically significant at the *p*-value smaller than 0.001.

**Summary.** In this section, we propose a new factor, the noisiness of context words, which influences the stability of medical concept embeddings. We use normalized entropy to estimate the proposed factor. The evaluation result shows a clear linear correlation between the proposed factor and the stability of medical concept embeddings. The correlations are consistent for all concepts, for low-frequency concepts, and for commonly used low-frequency concepts. This result supports the claim that when a concept has coherent context words, it is likely to result in high-quality embeddings, regardless of high or low frequency in a corpus. Moreover, the proposed factor provides an empirical reasoning of the surprising high stability of low-frequency concepts.

## 5 SUGGESTIONS ON IMPROVING THE STABILITY

We have shown the negative linear correlation in the evaluation. It indicates that decreasing the noisiness of context words improves the stability of medical concept embeddings. In this section, we make suggestions on ways to adjust the noisiness of context words, while maintaining the semantic meanings of concepts implied in a training corpus.

The first approach is to utilize hierarchical relations of medical concepts from a knowledge base. All medical concepts are linked to UMLS, and UMLS provides hierarchical relations among medical concepts. Hierarchical relations are a tree-like structure and connections indicate a parent-child (IsA) relation between concepts. In a parent-child relation, while the meaning of child concept is more specific than the meaning of parent concept, their meanings are greatly overlapped with each other. For example, Knee replacement (CUI: C0086511) is a child concept of Knee joint operation (CUI: C0187769) in UMLS, and the meaning is knee joint operation for replacement. The noisiness of context words of low-stability concepts can be adjusted using the context words of parent and/or child concepts.

Another approach to decreasing the noisiness of context words is segregating a training corpus into smaller document clusters. Context words tend to be more consistent in topic-specific document corpus and topic-specific sub-collection can be identified using document clustering techniques. Clustered documents provide information representing the importance of context words,

depending on which documents they appear. This information can be utilized to filter the context words.

## 6 GENERALIZATION

We have evaluated the proposed factor with concept embeddings trained with medical textual data. In this section, we conduct an experiment with embeddings trained on text from general domain. We use two datasets: Reuters-21578 (Reuters) and Wikipedia abstract (Wiki). Reuters consists of news articles under pre-defined topical categories. Wiki is much larger than Reuters and contains various topics.

In general domain, word embeddings are more prevalent than concept/entity embeddings. We train word embeddings without concept identification, and all trained word embeddings are subject to evaluations. Likewise, the skip-gram model with negative sampling is used. After the training, every word that has a frequency greater than minimum frequency threshold (set to 5) in a dataset has a word embedding.

Pearson correlation coefficients on Reuters and Wiki are $-0.412$ and $-0.495$, respectively. The proposed factor shows the moderate correlations with the stability of word embeddings in both datasets. The correlation values are statistically significant by $p < 0.001$. This result demonstrates that the linear relation of the proposed factor with the stability is not limited to medical-domain text, and it extends to general-domain text.

## 7 CONCLUSION

We analyze the stability of medical concept embeddings with respect to frequency. The analysis shows that low-frequency concepts can achieve high stability even though there is the limited number of training instances for them. Motivated by the surprising high stability of low-frequency concepts, we propose a new factor, the noisiness of context words, influencing the stability of medical concept embeddings. We measure the noisiness of context words using normalized entropy. The evaluation result shows a clear correlation between the proposed factor and the stability of medical concept embeddings. The result is consistent for all concepts, low-frequency concepts, and commonly used low-frequency concepts. This work helps understand how low-frequency concepts can result in high-stability embeddings like high-frequency concepts do.

## REFERENCES

[1] Edgar Altszyler, Mariano Sigman, and Diego Fernández Slezak. 2018. Corpus Specificity in LSA and Word2vec: The Role of Out-of-Domain Documents. In *Proceedings of The Third Workshop on Representation Learning for NLP*. Association for Computational Linguistics, 1–10.

[2] Maria Antoniak and David Mimno. 2018. Evaluating the stability of embedding-based word similarities. *TACL* 6 (2018), 107–119.

[3] Alan R Aronson. 2001. Effective mapping of biomedical text to the UMLS Metathesaurus: the MetaMap program.. In *Proceedings of the AMIA Symposium*. American Medical Informatics Association, 17.

[4] Andrew L Beam, Benjamin Kompa, Inbar Fried, Nathan P Palmer, Xu Shi, Tianxi Cai, and Isaac S Kohane. 2018. Clinical Concept Embeddings Learned from Massive Sources of Medical Data. *arXiv preprint arXiv:1804.01486* (2018).

[5] Piotr Bojanowski, Edouard Grave, Armand Joulin, and Tomas Mikolov. 2017. Enriching Word Vectors with Subword Information. *Transactions of the Association for Computational Linguistics* 5 (2017), 135–146.

[6] Xiangrui Cai, Jinyang Gao, Kee Yuan Ngiam, Beng Chin Ooi, Ying Zhang, and Xiaojie Yuan. 2018. Medical Concept Embedding with Time-Aware Attention. In *IJCAI*. 3984–3990.

[7] Hugo Caselles-Dupré, Florian Lesaint, and Jimena Royo-Letelier. 2018. Word2vec applied to recommendation: hyperparameters matter. In *RecSys*. 352–356.

[8] Billy Chiu, Gamal Crichton, Anna Korhonen, and Sampo Pyysalo. 2016. How to train good word embeddings for biomedical NLP. In *Proceedings of the 15th Workshop on Biomedical Natural Language Processing*. 166–174.

[9] Edward Choi, Mohammad Taha Bahadori, Elizabeth Searles, Catherine Coffey, Michael Thompson, James Bost, Javier Tejedor-Sojo, and Jimeng Sun. 2016. Multi-layer representation learning for medical concepts. In *KDD*. ACM, 1495–1504.

[10] Edward Choi, Mohammad Taha Bahadori, Le Song, Walter F Stewart, and Jimeng Sun. 2017. GRAM: graph-based attention model for healthcare representation learning. In *KDD*. ACM, 787–795.

[11] Youngduck Choi, Chill Yi-I Chiu, and David Sontag. 2016. Learning low-dimensional representations of medical concepts. *AMIA Summits on Translational Science Proceedings* 2016 (2016), 41.

[12] Lance De Vine, Guido Zuccon, Bevan Koopman, Laurianne Sitbon, and Peter Bruza. 2014. Medical semantic similarity with a neural language model. In *CIKM*. ACM, 1819–1822.

[13] Johannes Hellrich and Udo Hahn. 2016. Bad company-neighborhoods in neural embedding spaces considered harmful. In *Proceedings of COLING 2016, the 26th International Conference on Computational Linguistics: Technical Papers*. 2785–2796.

[14] Megan Leszczynski, Avner May, Jian Zhang, Sen Wu, Christopher R Aberger, and Christopher Ré. 2020. Understanding the Downstream Instability of Word Embeddings. *MLSys* (2020).

[15] Omer Levy, Yoav Goldberg, and Ido Dagan. 2015. Improving distributional similarity with lessons learned from word embeddings. *TACL* 3 (2015), 211–225.

[16] Nut Limsopatham, Craig Macdonald, and Iadh Ounis. 2014. Modelling Relevance towards Multiple Inclusion Criteria when Ranking Patients.. In *CIKM*. ACM, 1639–1648.

[17] Eneldo Loza Mencıa, Gerard de Melo, and Jinseok Nam. 2016. Medical concept embeddings via labeled background corpora. In *Proceedings of the 10th Language Resources and Evaluation Conference (LREC 2016), Paris, France*.

[18] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In *NIPS*. 3111–3119.

[19] Mark A Musen, Natalya F Noy, Nigam H Shah, Patricia L Whetzel, Christopher G Chute, Margaret-Anne Story, Barry Smith, and NCBO team. 2011. The national center for biomedical ontology. *Journal of the American Medical Informatics Association* 19, 2 (2011), 190–195.

[20] Denis Newman-Griffis, Albert M Lai, and Eric Fosler-Lussier. 2018. Jointly Embedding Entities and Text with Distant Supervision. In *Proceedings of The Third Workshop on Representation Learning for NLP*. 195–206.

[21] Jeffrey Pennington, Richard Socher, and Christopher Manning. 2014. Glove: Global vectors for word representation. In *EMNLP*. 1532–1543.

[22] Bénédicte Pierrejean and Ludovic Tanguy. 2019. Investigating the stability of concrete nouns in word embeddings. In *Proceedings of the 13th International Conference on Computational Semantics-Short Papers*. 65–70.

[23] Luca Soldaini and Nazli Goharian. 2016. Quickumls: a fast, unsupervised approach for medical concept extraction. In *MedIR Workshop, SIGIR*.

[24] Chih-Hsuan Wei, Hung-Yu Kao, and Zhiyong Lu. 2013. PubTator: a Web-based text mining tool for assisting Biocuration. *Nucleic Acids Research* 41 (07 2013).

[25] Laura Wendlandt, Jonathan K Kummerfeld, and Rada Mihalcea. 2018. Factors Influencing the Surprising Instability of Word Embeddings. In *NAACT-HLT*, Vol. 1. 2092–2102.

[26] Zhiguo Yu, Trevor Cohen, Byron Wallace, Elmer Bernstam, and Todd Johnson. 2016. Retrofitting word vectors of mesh terms to improve semantic similarity measures. In *Proceedings of the Seventh International Workshop on Health Text Mining and Information Analysis*. 43–51.

[27] Yongjun Zhu, Erjia Yan, and Fei Wang. 2017. Semantic relatedness and similarity of biomedical terms: examining the effects of recency, size, and section of biomedical publications on the performance of word2vec. *BMC medical informatics and decision making* 17, 1 (2017), 95.

## 8 APPENDIX

In the skip-gram model, different hyperparameters affect the quality of word embeddings. We evaluate the proposed factor on multiple sets of word embeddings that are trained with different hyperparameter values. We test five hyperparameters: window size (W), epoch (E), number of negative samples in negative sampling (N), smoothing parameter for negative sample distribution in negative sampling (M), and subsampling rate (S). Default values for W, E, N, M, and S are 7, 50, 5, 0.75, and 0.001, respectively. For each hyperparameter, two additional values are tested, while the rest

**Table 3: Average (standard deviation) of the stability of medical concept embeddings when different hyperparameter values are set in the training of word embeddings**

| Hyper-parameter | W:E:N:M:S | Avg (stdev) of stability | |
|---|---|---|---|
| | | cui2vec | NLM |
| *Default* | *7:50:5:0.75:0.001* | *0.595 (0.208)* | *0.627 (0.181)* |
| Window size (W) | 5 | 0.595 (0.205) | 0.636 (0.177) |
| | 10 | 0.589 (0.210) | 0.611 (0.187) |
| Epoch (E) | 30 | 0.564 (0.226) | 0.622 (0.189) |
| | 100 | 0.619 (0.197) | 0.622 (0.182) |
| Number of NS (N) | 10 | 0.618 (0.203) | 0.656 (0.177) |
| | 15 | 0.625 (0.202) | 0.669 (0.175) |
| Smoothing (M) | 0.0 | 0.579 (0.229) | 0.649 (0.179) |
| | 1.0 | 0.568 (0.209) | 0.611 (0.185) |
| Subsampling rate (S) | 0.01 | 0.604 (0.206) | 0.632 (0.181) |
| | 0.0001 | 0.579 (0.208) | 0.609 (0.18) |

hyperparameters are the default values. Likewise, we train medical concept embeddings using the two approaches, cui2vec and NLM.

Table 3 presents the average and the standard deviation of the stability. Marginal changes are observed in both cui2vec and NLM when different hyperparameter values are used.

Table 4 shows Pearson correlation coefficient of the proposed factor with the stability when embeddings are trained with different hyperparameter values. The correlation coefficients present small changes. This result makes sense that the stabilities are barely changed by the different hyperparameter values (Table 3)

Across the different hyperparameter values, the proposed factor shows a moderate linear correlation with the stability of medical concept embeddings. It is worth noting the correlations with the window size (W). The different values of W directly changes the noisiness of context words as the window size determines context words. When different window sizes are used in training of word embeddings, the proposed factor consistently shows the moderate correlations.

**Table 4: Pearson correlation coefficients of the noisiness of context words with the stability of medical concpet embeddings. Embeddings are trained by different hyperparameter values. The strongest correlation in cui2vec and NLM is indicated as \*. All correlation values are statistically significant at $p < 0.001$.**

| Hyper-parameter | W:E:N:M:S | Pearson coefficient | |
|---|---|---|---|
| | | cui2vec | NLM |
| *Default* | *7:50:5:0.75:0.001* | *-0.593* | *-0.477* |
| Window size (W) | 5 | -0.577 | -0.484 |
| | 10 | -0.605 | -0.470 |
| Epoch (E) | 30 | -0.618* | -0.478 |
| | 100 | -0.584 | -0.501 |
| Number of NS (N) | 10 | -0.564 | -0.471 |
| | 15 | -0.556 | -0.469 |
| Smoothing (M) | 0.0 | -0.519 | -0.335 |
| | 1.0 | -0.601 | -0.506* |
| Subsampling rate (S) | 0.01 | -0.590 | -0.484 |
| | 0.0001 | -0.578 | -0.454 |