

Homework #2

Instructor: Moontae Lee

Total Points: 170

Policy

1. HW2 is due by 11/30/2020 11:59PM in Central Time. One submission per each group.
2. You are allowed to work individually or as a group of up to three students.
3. Having wider discussions is not prohibited. Put all the names of students beyond your group members. However individual students/groups must write their own solutions.
4. Put your write-up and results from the coding questions in a single pdf file. Compress the source codes into a zip file. Each student/group must submit only two files. (You will lose the points if the answers for coding questions are not included the pdf report)
5. If you would include some graphs, be sure to include the source codes together that were used to generate those figures. Every result must be reproducible.
6. Maximally leverage Piazza to benefit other students by your questions and answers. Try to be updated by checking notifications in both Piazza and the class webpage.
7. Late submissions will be penalized 20% per each late day. As this is the last assignment in the semester, only one late day is allowed. For HW2, it will be 12/1/2020 11:59pm.

Problem 1: Word-vector Embedding

[25 points]

In class you have learned distributed representations which try to encode each word as a vector in a multi-dimensional Euclidean vector space. Answer for the following questions:

- (a) A naive encoding is to use $|V|$ -dimensional representation for each word where V is the set of vocabulary. As given in the lecture, you can construct such representations by counting either the co-existence (0-1) or the co-occurrence (frequency) with other words in the data. Explain the benefits and limitations of these approaches.
- (b) Suppose another encoding provides you $v_{man} = (0.3, 0.1, 0.4)$ and $v_{woman} = (0.3, 0.1, -0.6)$. Compute the cosine similarity between *man* and *woman* as learned in the class. Interpret some of the three dimensions by comparing lexical semantics of these two words.
- (c) Recall two word-vectors given in (b). Assuming $v_{boy} = (-0.7, -0.9, 0.3)$, guess the best word-vector for *girl*. Try to come up with vector operations that can evaluate v_{girl} in terms of addition/subtraction(s) of v_{boy} , v_{man} , and v_{woman} , thereby justifying your guess.

Word2vec is the most popular word-vector embedding that brings up innovations for various applications in Natural Language Processing. The following questions ask basic understanding about word2vec's theoretical foundations.

- (d) Given a target word $t \in V$ and a context word $c \in V$, the skip-gram models the conditional probability $p(c|t)$ by the following formula:

$$p(c|t) = \frac{\exp(u_c^T v_t)}{\sum_{w \in V} \exp(u_w^T v_t)}.$$

Explain why the exponential function is necessary and how this formula properly converts relationships in the vector space into a probability distribution.

- (e) In class, we derive the partial derivative of the log-likelihood version of the above equation with respect to v_t . Evaluate the partial derivative of it with respect to u_c . In other words, compute $\frac{\partial \log p(c|t; u_c, v_t)}{\partial u_c}$. Then explain how to learn word-vector embeddings.

Problem 2: Part-Of-Speech Tagging and Parsing [25 points]

In order for POS tagging, supervised dataset is generally required. Each example in the data consists of a sentence instance and the true label: tags that mark the true POS for each word. Once models learn from the data, then it can predict the most-likely POS tags of each word in unseen examples. Penn Treebank is the most popular supervised dataset.

- (a) Machine learning models predict well the labels of instances if they are already seen during the training process. Construct a simple sentence which is a part of supervised dataset, but the learned POS tagger could incorrectly predict its POS tags when testing on the same sentence.
- (b) All word *to* in the Penn Treebank is tagged simply as *TO* rather than as a precise POS. Explain potential problems by making your own examples that includes *to*. If you try to make an elaborated POS tagger that can distinguish different syntactic roles of *to* in your examples, what could you do?

Answer the following questions about parsing given the 9 rules: 1) $S \rightarrow NP VP$; 2) $S \rightarrow VP$; 3) $NP \rightarrow Det NP$; 4) $NP \rightarrow Proper-Noun Noun$; 5) $VP \rightarrow Verb NP$; 6) $Det \rightarrow the$; 7) $Noun \rightarrow run \mid marathon$; 8) $Verb \rightarrow run$; 9) $Proper-Noun \rightarrow Chicago$.

- (c) Show a possible bottom-up parsing for the sentence: *"Run the Chicago marathon"*.
- (d) Show a possible top-down parsing for the sentence: *"Run the Chicago marathon"*.
- (e) If you only draw a parse-tree of the sample sentence used in (c) and (d), can you tell which derivation algorithm you used between the top-down and the bottom-up approaches?

Problem 3: Programming Project

[100+20 points]

Word Sense Disambiguation (WSD) is a task to find the correct meaning of a word given context, which can be a building block for various high-level NLP tasks. As many words in languages have more than a single meaning, humans perform WSD with respect to various verbal and non-verbal signals. In this problem, you are going to implement a WSD system by using two different models: *ontological* model and *supervised* model. To start, read the English Lexical Sample Task written by Mihalcea, Chklovski and Kilgarrriff in the following link.¹

The data files are lightly preprocessed for the class project. They consist of training, validation, and test data provided with a XML formatted dictionary that describes commonly used senses for each word. Every lexical element in the dictionary contains multiple sense items, assigning one integer id per each sense. Briefly see the following example from our XML dictionary

```
<lexelt item="future.n" num="4">
  <sense id="1" wordnet="1" gloss="time to come" examples="In the future we will drive flying cars.
  | What will you do in the future." />
  <sense id="2" wordnet="2" gloss="verb tense" examples="The paper was written in future tense." />
  <sense id="3" wordnet="3" gloss="commodities" examples="He made his living trading in futures." />
  <sense id="4" wordnet="" gloss="personal time to come" examples="My future is bright. | What will
  you do with your future?" />
</lexelt>
```

It describes one lexical element: *future* (part-of-speech is noun) with its four different senses. Each sense has its own *gloss* (definition) and *examples* that are separated by | symbol. Each sense is also associated with the corresponding senses of WordNet 2.1. As our sense divisions and WordNet's are not identical, some of our senses could be mapped to multiple WordNet senses or possibly nothing (e.g., See the fourth sense item in the above example). Since the current NLTK is using WordNet 3, the sense mapping from 2.1 to 3.x could be useful.²

The training data specifies the correct sense of the target word providing its verbal context surrounding the target word. Each line of training data has the the following format:

```
word.pos | sense-id | prev-context %% target %% next-context
```

- **word** is the original form of the target word for which we are to predict the sense. You will use it to lookup the XML dictionary.
- **pos** is the POS where 'n', 'v', and 'a' stand for noun, verb, and adjective, respectively.
- **sense-id** is the integer number for the correct sense id defined in our dictionary.
- **prev-context** is the text given earlier than each of the target word occurrence.
- **target** is the actual occurrence of the target word. Note that the word "begin.v" could occur as "beginning" instead of "begin" to denote a participle at the given position.
- **next-context** is the text given later than each of the target word occurrence.

Note that sense-ids in the test data are all erased to 0 as those are what you should predict.

¹<http://citeseerx.ist.psu.edu/viewdoc/summary?doi=10.1.1.153.8426>

²<https://stackoverflow.com/questions/46950379/how-to-fetch-a-specific-version-of-wordnet-when-doing>

1 Ontological (dictionary-based) WSD

Dictionary-based approaches utilize definitions given in the dictionary. See the following example that tries to disambiguate "*pine cone*".

- pine (the context)
 1. a kind of **evergreen tree** with needle-shaped leaves
 2. to waste away through sorrow or illness
- cone (the target word)
 1. A solid body which narrows to a point
 2. Something of this shape, whether solid or hollow
 3. Fruit of certain **evergreen trees**

As bold faced in the above, 3rd sense of the target word matches the most with the 1st sense of the context word among all possible combinations. This process shows the original Lesk algorithm to disambiguate senses based only on the cross-comparing the definitions. However, rich examples given in the dictionary can be utilized to extend this model for better matching.

1. Design a metric that rewards consecutive overlaps more. One overlap of two consecutive words must get higher scores than two distant overlaps of a single word. Note that there would be morphological variations in the definitions and examples. To increase the matching, stemming or lemmatizing could be useful.³
2. Implement a dictionary-based WSD system that disambiguates the sense by comparing the definitions of the target word to the definitions of relevant words in the context. Your design decision of choosing relevant words will determine the performance of the dictionary-based system in combining with the metric you designed above.
3. Because we mainly use glosses and examples in the dictionary to figure out the correct senses, no training is necessary for the Simple Lesk WSD. If you want to try the Corpus Lesk WSD (extension), try to augment the dictionary by the training data. If you think that those are not enough to achieve competitive accuracy, feel free to use the WordNet dictionary to further improve the performance.⁴
4. If no training process is involved, you could verify the performance of your Simple Lesk WSD system on the entire training set. If you want to compare the performance of various WSD systems like your Corpus Lesk or supervised WSD in the next section, test on the same validation set that we provide. You should also submit prediction results on the test data for every model that you would try.

³You can find relevant tools: stemmer and lemmatizer in NLTK and WordNet.

⁴If you would use WordNet, be careful in the version difference as stated in the introduction.

2 Supervised WSD

This section describes a simple probabilistic approach called the Naive-Bayes model. The model takes a word in context as an input and outputs a probability distribution over predefined senses, indicating how likely each sense would be the correct meaning of the target word within the given context. Specifically, it picks the best sense by the following equation:

$$\hat{s} = \operatorname{argmax}_{s \in S(w)} p(s|\vec{f})$$

In the above equation, $S(w)$ is the predefined set of senses for the target word w and \vec{f} is a feature vector extracted from the context surrounding w . Thus the equation says that we are going to choose the most probable sense as the correct meaning of w . By Bayes rule,

$$p(s|\vec{f}) = \frac{p(\vec{f}|s)p(s)}{p(\vec{f})}.$$

As the denominator does not change with respect to $s \in S(w)$, the best sense \hat{s} is given by

$$\hat{s} = \operatorname{argmax}_{s \in S(w)} p(s|\vec{f}) = \operatorname{argmax}_{s \in S(w)} p(\vec{f}|s)p(s).$$

Here the model *naively* assumes⁵ that each feature in the feature vector \vec{f} is conditionally independent given the sense of the word s . The assumption allows us to evaluate $p(\vec{f}|s)$ by

$$p(\vec{f}|s) = \prod_{j=1}^n p(f_j|s) \quad \text{where } f = (f_1, f_2, \dots, f_n).$$

In other words, the probability of a feature vector given sense can be estimated by the product of the probability of its individual features given that sense under our assumption. Hence,

$$\hat{s} = \operatorname{argmax}_{s \in S(w)} p(\vec{f}|s)p(s) = \operatorname{argmax}_{s \in S(w)} p(s) \prod_{j=1}^n p(f_j|s)$$

What you have to implement for this model is given in the following instructions.

1. To train the above model, you should learn the model parameters: 1) the prior probability of each sense $p(s)$ and 2) the individual feature probabilities $p(f_j|s)$. Those are computed by the Maximum Likelihood Estimation (MLE) which purely counts the number of actual occurrences in the training set. Particularly for the i -th sense s_i of a word w ,

$$P(s_i) = \frac{\text{count}(s_i, w)}{\text{count}(w)} \quad P(f_j|s_i) = \frac{\text{count}(f_j, s_i)}{\text{count}(s_i)}$$

For instance, assume there are 1,000 training examples corresponding to the word “bank”. Among them, 750 occurrences stand for $bank_1$ which covers the financial sense, and 250 occurrences for $bank_2$ which covers the river sense. Then the prior probabilities are

⁵This is why the model is called Naive-Bayes.

$$P(s_1) = \frac{750}{1000} = 0.75 \quad P(s_2) = \frac{250}{1000} = 0.25$$

If the first feature “credit” occurs 195 times within the context of $bank_1$, but only 5 times within the context of $bank_2$,

$$P(f_1 = \text{“credit”} | s_1) = \frac{195}{750} = 0.26 \quad P(f_1 = \text{“credit”} | s_2) = \frac{5}{250} = 0.02$$

2. The performance of your WSD system would rely more on how to generate feature vectors from the context. Note that target words are always provided within sufficiently long sentence(s). As the above example shows, extracting informative words from surrounding context allows the model parameters to discriminate unlikely senses from the correct sense. In our model, this process of deciding model parameters becomes *training*. **You have to train a separate model per each target word in the training data.**
3. When initially training your model, make sure that you never use the validation/test data. Note that the correct sense-ids given in the test data are deliberately erased to 0, which means those are no longer true labels. Instead of marking the predicted senses directly on the testing file, you must generate a separate output file consisting only of the predicted sense-ids, one id per line in each test data.
4. To achieve quality performance, smoothing is necessary. Implement either add-1 or add- λ smoothing. If you want to compare the performance of multiple models (e.g., different λ 's), feel free to use a validation set, which is randomly reserved from the original test data. Since the true senses are alive in the validation data, testing on the validation set will let you guess the true performance on the unseen data. Note that you must not train on the validation set if you want to validate the performance. However, you can add the validation set to the training data for predicting the best senses of the test data.

3 Scoring and Extensions

We use *accuracy*⁶ as a score. Since each of possible senses is already specified by a different sense-id, and no examples has multiple senses at the same time, a single prediction will be counted as incorrect one unless the prediction is equivalent to the ground-truth sense tag.

1. Assuming the given word in a test example has k different senses based on our XML dictionary, the prediction file must consist of a $1 - k$ integer number per line. Concretely, if the test set consists of three examples where each example has 7, 3, and 5 different senses, your system should output one line for each of three test examples like the following.

7
1
4

⁶Accuracy = # of correct predictions / # of total predictions

2. (+5 pts) Implement the Corpus Lesk algorithm by augmenting the dictionary with the training data. Report your improved performance against the Simple Lesk algorithm.
3. (+5 pts) Instead of hard-comparing to a single correct sense, you could design soft-scoring scheme that partially votes to each sense with respect to its confidence based on your model. For the supervised WSD using the Naive-Bayes model, it is easy to vote partially because the system guesses the correctness of each sense as a probability distribution, whereas you may have to do some normalization for soft-scoring in the dictionary-based method.⁷ Note that this scoring is an expected score: for example, if your best answers are *sense-1* with 70% confidence and *sense-2* with 30% confidence, you gain only 0.7 (rather than 1.0) if *sense-1* is a right answer, whereas you gain only 0.3 (rather than 1.0) if *sense-2* is a right answer. Evaluate the prediction result **on the validation set** and compute the average accuracy. Analyze the difference between the two scoring schemes and discuss which one seems more beneficial with supporting reasons.
4. (+10 pts) Use embeddings via Spacy package in Python. Install Spacy and download the pre-trained embedding models by “python -m spacy download en_core_web_md”. Then you can retrieve individual word-vectors given a sentence or its sentence-vector in terms of the average of the word vectors. Improve your ontological WSD and supervised WSD by incorporating these word-vector information. Feel free to use the following script for this open-ended extension.

```

1 import spacy
2 # Load the spacy model that you have installed.
3 model = spacy.load('en_core_web_md')
4 # Process a sentence given the pre-trained model.
5 embeddings = model("You are working on the second homework.")
6 # Extract a word-vector for the 7-th word homework.
7 embeddings[6].vector
8 # Get a sentence-vector as a mean of the individual word vectors.
9 embeddings.vector

```

4 What to submit?

Minimally you should implement the Simple Lesk algorithm for the ontological WSD and the Naive-Bayes algorithm with add-1 smoothing for the supervised WSD. After experiments, write a short PDF report (max 4 pages) that consists of the followings in addition to **your codes and predictions on the test data**. Add accordingly if you do some extensions. (List any software that you did not write by yourself. Note that using any pre-built WSD is not allowed)

- (a) Explain all WSD systems that you have built. Ideally two systems: the Simple Lesk and the Corpus Lesk for ontological WSD. Another two systems: add-1 and add- λ smoothing Naive-Bayes for supervised WSD.

⁷Recall 1-(d) that explains how to normalize the score, getting a probability distribution.

- (b) Try various scoring functions and different feature engineering. Pick the best one for each WSD system, providing several intuitive real examples chosen from the training data that can justify your design decisions.
- (c) Report the comparative performance among your ontological WSD systems with table/-graphs by testing on the validation set. Report the comparative performance similarly among your supervised WSD systems. Note that there must be a baseline WSD system that always predicts to the most frequent sense. No comparisons with the baseline cannot justify performance of your systems. Finally compare the entire WSD systems, reporting clearly labeled tables/graphs with a written summary of the results.
- (d) Include observations that you achieve during the experiment. One essential discussion is to analyze informative features based on the real examples. In addition, Discuss the difference between the supervised and the dictionary-based WSD systems. Which system is more appropriate for which cases based on the real examples chosen from the data.
- (e) Report your additional findings if you decide to implement some of the extensions.