

# NLPHW3

## Word Sense Disambiguation (WSD)

Manoochehr Joodi Bigdello - 1860273

### Project Overview

Word sense disambiguation (WSD) is the ability to identify the meaning of words in context. WSD is an important topic in NLP because most of the human languages are ambiguous, so that many words can be interpreted in multiple ways depending upon the context where they appear and it's a challenge for machines to understand the meanings (senses) behind the words. This task is Training multi-inventory WSD system capable of performing both fine-grained and coarse-grained WSD. Dataset is Raganato format sense-tagged English corpora (semcor and semcor+omsti) with mappings from WordNet to BabelNet synsets and BabelNet to LexNames and WordNetDomains.

### Preprocessing

In preprocessing phase, after parsing xml files and perform all mappings, getting data from WordNet and converting all data to id (we have dictionary for RawWords, SensesWords, POS, Lex, WnDomain), we build 4 different masks to filter our results of logits that we are going to feed in "sparse\_softmax" layer and argmax for prediction.

- Pos2Sense mask: it's a (size of POS dictionary)\*(size of SensesWords dictionary) to mask every possible sense to its POS tag.
- Lex2Sense mask: it's a (size of Lex dictionary)\*(size of SensesWords dictionary) to mask every possible sense to its LexName tag.
- WnDomain2Sense mask: it's a (size of WnDomain dictionary)\*(size of SensesWords dictionary) to mask every possible sense to its WnDomain tag.
- lemma2Sense mask: this is a little bit different from other masks, because of Tensorflow limits in GPU we have to perform this mask inside batches so it's a 3D mask (size of batch)\*(size of sentences with padding)\*(size of SensesWords dictionary) to mask every possible senses of a word to its lemma, this helps us to perform final evaluation as argmax of only the candidate synsets for the lemma we are trying to disambiguate. ( E.g. car#noun has C=[bn:00007309n, bn:00015785n, bn:00015786n, bn:00015787n, bn:00014462n] then synset = argmax(prob[C]))

### Models

We address this problem using of neural architecture bidirectional Long Short Term Memory (LSTM). We used all-words model which rely on sequence learning. We also used POS, LexName and WnDomain tags to improve the performance. As mentioned in homework assignment file we implement the MFS as our Back Off Strategy in case of no prediction for instances.

As embedding layer we used sense embedding's from our NLP HW 2 assignment and also embedding layer of Tensorflow for understanding the impact of sense embedding's on WSD.

We used BD-LSTM layer and then concat all outputs of forward and backward LSTM and feed that as input for other BD-LSTM layer on top of this one.

We also put attention layer + hidden layer on top of our sense BD-LSTM for getting information of all words for better prediction (as mentioned in paper).

For possible performance improvement we also used hierarchical modeling. We used POS, Lex and WnDomain tags as masking for results of sense layer and then feed the masked results to sparse\_softmax, we used sparse\_softmax because each lemma word can have exactly one sense. For minimizing we used gradient based minimizing with clip and our total\_loss to minimize was the sum of all sparse\_softmax parts.

### **We implement 7 different models for comparison:**

- 1- **MFS Tagger**: as base line we get MFS synset of lemma from WordNet.
- 2- **Basic Tagger**: we just used an Embedding layer + BD-LSTM layer for Sense + attention layer
- 3- **Multitask<sub>POS</sub> Tagger**: we used Embedding layer + BD-LSTM layer for POS (we compute loss for this layer with sparse\_softmax as loss\_pos) + BD-LSTM layer for Sense (input of BD-LSTM is output of POS BD-LSTM) with attention and hidden layer on top + POS masking (we use Pos2Sense mask to filter the possible senses based on POS) + sparse\_softmax to compute loss\_sense and predictions. We minimize the sparse\_softmax cross entropy loss and also regularize the weights using L2 regularization. We also clipped the gradients.
- 4- **Multitask<sub>Lex</sub> Tagger**: same as Multitask<sub>POS</sub> we just use LexNames instead of POS with Lex2Sense mask.
- 5- **Multitask<sub>WnDomain</sub> Tagger**: same as Multitask<sub>POS</sub> we just use WnDomains instead of POS with WnDomain2Sense mask.
- 6- **Multitask<sub>POS+Lex</sub> Tagger**: we used Embedding layer + BD-LSTM layer for POS (we compute loss for this layer with sparse\_softmax as loss\_pos) + BD-LSTM layer for LexNames (we compute loss for this layer with sparse\_softmax as loss\_lex) + BD-LSTM layer for Sense (input of BD-LSTM is output of POS BD-LSTM) with attention and hidden layer on top + POS and LexName masking (we use Pos2Sense and Lex2Sense mask multiplication to filter the possible senses based on POS and LexNames) + sparse\_softmax to compute loss\_sense and predictions. We minimize the sparse\_softmax cross entropy loss and also regularize the weights using L2 regularization. We also clipped the gradients
- 7- **Multitask<sub>POS+WnDomain</sub> Tagger**: same as **Multitask<sub>POS+Lex</sub>** we just use WnDomains instead of LexNames and WnDomain2Sense mask instead of Lex2Sense mask.

## Training

In training we used different hyperparameters as shown in table 1, then we choose the best one, for tuning we used 5000 sentence of SemCor dataset and semeval2007 as development set. For optimizer we decided to use the AdadeltaOptimizer with learning rate of 1 as mentioned in reference paper of optimizer. We trained every model in 10 Epoch, we just used the SemCor dataset (37K sentence) for training with 2 mode, first with Tensorflow Embedding Layer and then with Sense Embedding's of NLPHW2 Assignment.

## Test phase

We test our models and configurations on all 5 test set, Senseval-2 all-words, Senseval-3 all-words, SemEval 2007, SemEval 2013, and SemEval 2015, we also tested on ALL of them together to have general view of our proposed models. We predict our results in BabelNet format and then used mapping files to get LexNames and WnDomains. We report the results in 3 different cases BabelNetId's, LexNames and WnDomains output format, we also report the results on both with sense embedding's and without sense embedding's. The results were shown in Tables 2-4, bold fonts for comparing between with sense embedding's and without sense embedding's, and red highlighting is for comparing different models.

## Conclusion

As we trained our model with different configurations, we can see that the **Multitask<sub>POS+WnDomain</sub>** model is better than other models and we think it's because in training dataset after mapping the WnDomains are more accurate than LexName, also because of masking possible synsets for candidate lemma makes this model better than **Multitask<sub>POS</sub>**, **Multitask<sub>WnDomain</sub>** and **Multitask<sub>Lex</sub>**. Using of sense embedding's in most of the cases increase the performance but it's not too much and we think that it can be because of our poor model in HW2.

## References and Link

- 1- Neural Sequence Learning Models for Word Sense Disambiguation, Raganato, Delli Bovi and Navigli, EMNLP 2017
- 2- Tensorflow.com

## Charts and Tables

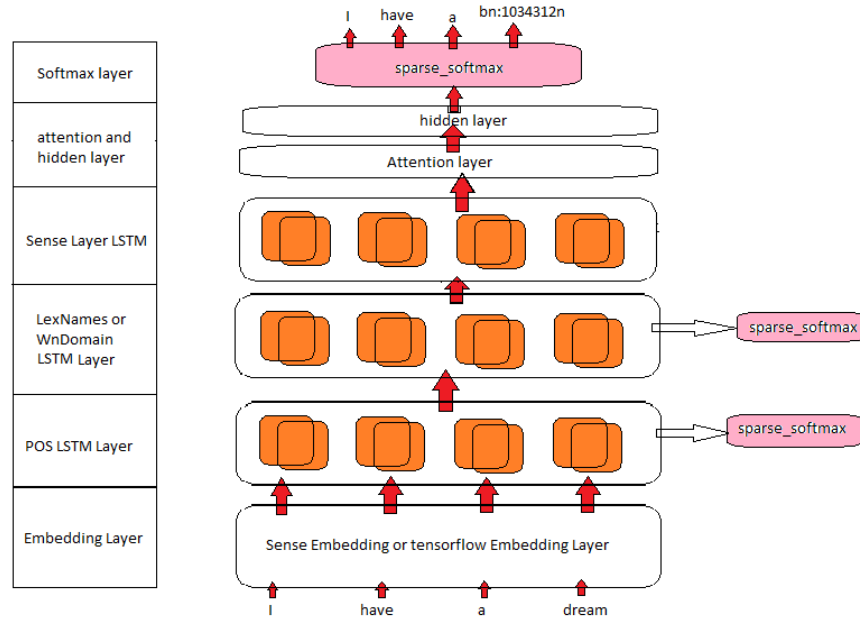


Figure 1. Our General model structure

Table 1. Hyper parameter tuning values

Hyperparametr	values	Best ones
Batch Size	4, 8, 16	4
CLIP	5, 10	10
Embedding dimension	100	100
Number of LSTM Layers	1, 2	2
Hidden Layers	128, 256	256

Table 2. F1 accuracy of different models on each test dataset with and without sense embedding in **BabelNet** format

	Senseval2		Senseval3		SemEval2007	
	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding
MFS Tagger	-	49.2	-	47.9	-	37.1
Basic Tagger	63.5	<b>64.1</b>	<b>61.9</b>	<b>62.2</b>	48.3	<b>49.8</b>
Multitask <sub>POS</sub> Tagger	<b>59.1</b>	57.9	<b>60.9</b>	60.2	<b>49.4</b>	47.4
Multitask <sub>Lex</sub> Tagger	<b>54.2</b>	53.9	<b>54.5</b>	54.4	<b>48.7</b>	48.1
Multitask <sub>WnDomain</sub> Tagger	<b>58.3</b>	58.1	<b>60.3</b>	58.3	<b>49.6</b>	48.7
Multitask <sub>POS+Lex</sub> Tagger		49		50.6		47.9
Multitask <sub>POS+ WnDomain</sub> in Tagger	<b>63.7</b>	63	60.7	<b>60.8</b>	48.7	<b>50.7</b>

	SemEval2013		SemEval2015		ALL	
	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding
MFS Tagger	-	<b>62.7</b>	-	44.3	-	<b>50.5</b>
Basic Tagger	59.2	<b>60.9</b>	59.7	<b>60</b>	60.6	<b>61.4</b>
Multitask <sub>POS</sub> Tagger	<b>58.3</b>	57.8	<b>58.2</b>	57.7	<b>60.6</b>	59.8
Multitask <sub>Lex</sub> Tagger	<b>58.8</b>	58.6	51.1	<b>51.2</b>	<b>54.7</b>	54.4
Multitask <sub>WnDomain</sub> Tagger	<b>53.3</b>	51.4	<b>52.4</b>	50.8	<b>56.3</b>	55
Multitask <sub>POS+Lex</sub> Tagger		53.2		46.5		50
Multitask <sub>POS+WnDomain</sub> Tagger	<b>60</b>	58.6	<b>59.3</b>	58.7	<b>60.6</b>	60

Table 3. F1 accuracy of different models on each test dataset with and without sense embedding in **LexNames** format

	Senseval2		Senseval3		SemEval2007	
	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding
MFS Tagger	-	58.7	-	55.5	-	43.7
Basic Tagger	81.1	<b>81.5</b>	77.3	<b>77.6</b>	67.2	<b>68.3</b>
Multitask <sub>POS</sub> Tagger	<b>76.7</b>	75.8	<b>76.5</b>	76	<b>68.3</b>	67.7
Multitask <sub>Lex</sub> Tagger	<b>74.2</b>	73.4	72.1	<b>72.4</b>	<b>67.9</b>	67.9
Multitask <sub>WnDomain</sub> Tagger	75.6	<b>76.3</b>	<b>76</b>	73.6	66.8	<b>67.4</b>
Multitask <sub>POS+Lex</sub> Tagger	69.1	<b>69.5</b>	<b>70.8</b>	70	65.8	<b>67.2</b>
Multitask <sub>POS+WnDomain</sub> Tagger	<b>81.2</b>	80.9	<b>76.3</b>	75.9	68.3	<b>69.2</b>

	SemEval2013		SemEval2015		ALL	
	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding
MFS Tagger	-	<b>74.9</b>	-	54.5	-	60
Basic Tagger	70.3	<b>72.3</b>	77.4	<b>78.1</b>	76.3	<b>77.1</b>
Multitask <sub>POS</sub> Tagger	<b>70.2</b>	69.8	<b>76.9</b>	76.4	<b>76.7</b>	76.2
Multitask <sub>Lex</sub> Tagger	<b>70.8</b>	70.5	<b>71.9</b>	71.7	<b>72</b>	71.9

<b>Multitask<sub>WnDomain</sub> Tagger</b>	<b>64.5</b>	62.9	<b>70.5</b>	70.15	<b>71.9</b>	71.1
<b>Multitask<sub>POS+Lex</sub> Tagger</b>		64.1		66.5		67.9
<b>Multitask<sub>POS+ WnDomain</sub> Tagger</b>	<b>71.6</b>	70.7	<b>77.7</b>	<b>78.4</b>	<b>76.6</b>	76.2

Table 4. F1 accuracy of different models on each test dataset with and without sense embedding in **WnDomains** format

	Senseval2		Senseval3		SemEval2007	
	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding
<b>MFS Tagger</b>	-	53	-	57.4	-	69.4
<b>Basic Tagger</b>	<b>58.8</b>	<b>59</b>	67.2	<b>67.5</b>	83.9	<b>84.1</b>
<b>Multitask<sub>POS</sub> Tagger</b>	<b>55.3</b>	54.2	<b>66.4</b>	65.8	84.1	<b>84.6</b>
<b>Multitask<sub>Lex</sub> Tagger</b>	57.9	<b>58.3</b>	<b>67.7</b>	67	84.5	<b>85</b>
<b>Multitask<sub>WnDomain</sub> Tagger</b>	54.7	54.7	<b>65.3</b>	63.9	84.3	<b>84.6</b>
<b>Multitask<sub>POS+Lex</sub> Tagger</b>		55.1		63.4		81.5
<b>Multitask<sub>POS+ WnDomain</sub> Tagger</b>	<b>58.7</b>	58.4	<b>65.7</b>	65.5	<b>85.2</b>	<b>85.5</b>

	SemEval2013		SemEval2015		ALL	
	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding	With sense embedding	Without sense embedding
<b>MFS Tagger</b>	-	<b>74.8</b>	-	49.7	-	59.6
<b>Basic Tagger</b>	<b>72.2</b>	<b>73.9</b>	57.1	<b>57.5</b>	65.3	<b>65.9</b>
<b>Multitask<sub>POS</sub> Tagger</b>	<b>71.5</b>	70	<b>56.6</b>	56.4	<b>65.5</b>	64.6
<b>Multitask<sub>Lex</sub> Tagger</b>	<b>71.1</b>	70.7	56.3	<b>56.4</b>	<b>65.1</b>	64.7
<b>Multitask<sub>WnDomain</sub> Tagger</b>	64.6	<b>64.7</b>	<b>52.6</b>	51	<b>61.2</b>	60.7
<b>Multitask<sub>POS+Lex</sub> Tagger</b>		66		50.7		60.7
<b>Multitask<sub>POS+ WnDomain</sub> Tagger</b>	71	<b>71.7</b>	<b>56.8</b>	56.6	<b>64.7</b>	64.7