# SEMI-SUPERVISED LABEL PROPAGATION FOR IDENTIFICATION OF DERIVATIONAL NOUNS IN ENGLISH

**Submitted by:**

Satanik Mitra (17IM91R04)

Samridha Kumar (16CS71P03)

**Project Mentor:**

Amrith Krishna

**Faculty:**

Prof. Niloy Ganguly

# ABSTRACT

Derivational nouns are widely used in the English language and is a prevalent means of productivity in the language. Derivational nouns are derived by affixing suffixes to the base form of verbs. For example, affixing *-er* to the base form of the verb *learn* results in the noun *learner*, and affixing *-or* to the base form of the verb *interrogate* results in the noun *interrogator.*

However, affixes that modify the mere morphological or syntactical role of a word in its usage are not considered derivational, rather inflectional.

Identifying derived words from a corpus is challenging. Usage of pattern matching approaches in strings are often inept for the task. Tasks that rely on string matching approaches alone, often result in a large number of false positives. For example, while the word 'postal' is generated from 'post', the word 'canal' is not generated from 'can'. String matching approaches often result in low recall as well, due to the variations in patterns in the derived and source word pairs, even for the same affix. Both 'postal' and 'minimal' are derived using the affix 'al', but the source word for postal is 'post', while the source word for minimal is 'minimum'.

In our work, we propose an approach for the identification and analysis of derivational nouns in the English Corpora. We obtain the rules for generation of derivational nouns and combine the diverse features ranging from them with the word embeddings for the candidate words. The lack of labelled data and other resources have led us to propose a graph based semi-supervised label propagation approach, using Modified Adsorption (MAD), to identify the usage of derivational nouns in a corpus and map them to their corresponding source words. The label propagation can be summazsed as:

1.  The nodes of the graph represent the labelled/unlabelled data points, and the weight of the edges is the measure of closeness or proximity with the other node.

2. Lables of a node propagate to all other nodes through edges.

3. Larger edge weights allow labels to travel through easier.

By leveraging on the rules obtained, we not only find the different pattern differences between the source and the derived word pairs, but also group pattern that are likely to emerge from the same affixes.

# CONTENTS

# OBJECTIVE

The entire work done can be broadly categorized into two objective goals. Our work has focused on achieving these objective goals.

1. Find the different pattern differences between source and derived word pairs, and group patterns from same affixes.
2. Assessment of MODIFIED ADSORPTION (MAD) algorithm with respect to prediction of derivational affixes.

The basic challenges faced in the process of achieving the objectives were:
1. Detection of False Positives (FP):

   There are some word pairs which though bear semblance to the derivational affix pattern, but are actually not *'derivational'*.

   e.g. *'Canal'* is not derived from the root *'Can'* affixing the morpheme *'-al'*.

   Our model was so trained such that these word pairs were identified as False Positives.

2. Correcting Low Recall:

   There are some derivational words which may bear the same affix but the source word is not similar.

   e.g. *'Post-al'* is derived from *'Post'*, but *'Minim-al'* is derived from *'Minimum'*.

   Our work has focused on detecting these word pairs which give a Low Recall to the trained model, and matching them to their corresponding correct source word.
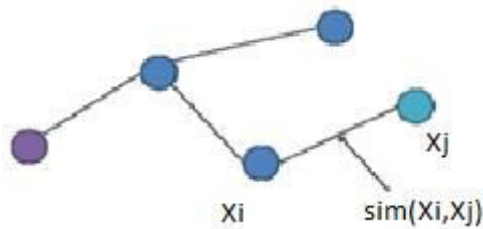
# DESCRIPTION

For the purpose of our work, we have adopted a Graph based Semi Supervised Learning (SSL) Approach for transductive learning.

Semi-supervised Learning was chosen because of the lack of labeled data. Semi-supervised algorithms learn from limited amounts of labeled data and vast amounts of freely available unannotated data. SSL also helps in generation of a more accurate decision boundary.

In a Graph based SSL, each example is associated with a vertex in an undirected graph and a weighted edge between two vertices represents similarity between the two corresponding examples. Graph-based semi-supervised algorithms often propagate the label information from the few labeled vertices to the entire graph. Most of the algorithms trade-off between *accuracy* (initially labeled nodes should retain those labels, relaxations allowed by some methods) with *smoothness* (adjacent vertices in the graph should be assigned similar labels). In a transductive setting, the algorithms output label information to the unlabeled data.

Among the various Graph based SSL algorithms present, we have chosen **Modified Adsorption** for the purpose of our work. It performs the multiclass classification required for our task, and can be parallelized and scaled to handle large data.

For our work, we have used the *Junto Label Propagation Toolkit* which contains an implementation of Modified Adsorption.



There are two stages in the complete Modified Adsorption algorithm:
1. Graph Construction
2. Label Inference

The next section describes the two stages and the entire methodology in detail.

# METHODOLOGY & EXPERIMENTS

This section describes in detail the entire methodology adopted for the purpose of our work. The work-flow can be divided into two categories:

1. Graph Construction
2. Label Inference

## Graph Construction

**STEP 1:** Extraction of source word – derived word pairs from the English Corpus.

A total of 982 source word – derived word pairs have been extracted from the English corpora. These word pairs consist of a mixture of derivational affixes, namely: *-ally, -ism, -ist, -ic, -y, -man, -like, -hood, -ation*. Each affix group consists of 98 word pairs, apart from the group corresponding to affixes *-ist* and *-y* which consist of 99 word pairs.

**STEP 2:** Representation of the word pairs as word vectors.

Each word of the word pair has been represented by a word vector of length 300. The word vectors for each word were derived using the *fastText* library. Facebook has released pre-trained fastText 300-dimensional word vectors for over 90 languages and our work makes use of that.



Fig. 2: Word Vectors for one source word – derived word pair.

**STEP 3:** Representation of a single word pair as a single node.

This step requires representing the entire word pair as a single word vector. For this purpose, we have calculated the average of the word vector for source word and derived word.

We can also take the word vector for the source word or the derived word as the word vector for the entire word pair.

**STEP 4:** Calculation of edge weights between the different nodes.

For the calculation of edge weights between the different nodes, we calculate the similarity, specifically *Cosine Similarity* between each vector i.e. each word pair (node).

For our work, we have considered edge weights between all unlabelled nodes and edge weights between each labelled node to the set of nodes belonging to same affix class. i.e. there is a edge between labeled nodes of -*ally* and the corresponding nodes of class -*ally,* and so on.

With the calculation of edge weights between the nodes, the construction of input graph is complete. The input graph constructed consists of 982 nodes and 20460 edges.

```
N978      N660      0.161595854261
N978      N661      0.22711829254
N978      N662      0.399691951378
N978      N663      0.260870120511
N978      N664      0.222198624115
N978      N665      0.175651648377
N978      N666      0.175319137669
N978      N667      0.171532768441
N978      N668      0.16849334859
N978      N669      0.285803187737
N978      N670      0.166318621902
N978      N671      0.214757153066
N978      N672      0.324437867168
N978      N673      0.147498841649
N978      N674      0.139333401321
N978      N675      0.251484073457
N978      N676      0.275148497771
N978      N677      0.288969440275
N978      N678      0.250831961909
N978      N679      0.336321682067
N978      N680      0.261044854178
N978      N681      0.317040636088
N979      N614      0.109469270998
N979      N615      0.241518588189
N979      N616      0.153463638514
N979      N617      0.147958097762
N979      N618      0.23342235266
N979      N619      0.0876443099189
N979      N620      0.156048233609
N979      N621      0.189962427648
N979      N622      0.289400124728
N979      N623      0.153392539853
```

Fig. 3: Nodes pairs with corresponding edge weights.

# Label Inference:

For the label inference task, the MAD algorithm was implemented on our constructed input graph. This was done using *Junto* toolkit, which provides implementation of the MAD algorithm.

The input graph constructed was given input to the *Junto* as *input_graph.*

A seed file was created, *seeds,* which consists of the labeled nodes. A *gold_labels* file was also created which was used to compare the predicted labels of the unlabeled nodes with the actual labels.

| N682 | L1 | 1.0 |
|------|------|-----|
| N683 | L1 | 1.0 |
| N684 | L1 | 1.0 |
| N685 | L1 | 1.0 |
| N686 | L1 | 1.0 |
| N687 | L1 | 1.0 |
| N688 | L1 | 1.0 |
| N689 | L1 | 1.0 |
| N690 | L1 | 1.0 |
| N691 | L1 | 1.0 |
| N692 | L1 | 1.0 |
| N693 | L1 | 1.0 |
| N694 | L1 | 1.0 |
| N695 | L1 | 1.0 |
| N696 | L1 | 1.0 |
| N697 | L1 | 1.0 |
| N698 | L1 | 1.0 |
| N699 | L1 | 1.0 |
| N700 | L1 | 1.0 |
| N701 | L1 | 1.0 |
| N702 | L1 | 1.0 |
| N703 | L1 | 1.0 |
| N704 | L1 | 1.0 |
| N705 | L1 | 1.0 |
| N706 | L1 | 1.0 |
| N707 | L1 | 1.0 |
| N708 | L1 | 1.0 |
| N709 | L1 | 1.0 |
| N710 | L1 | 1.0 |
| N711 | L1 | 1.0 |
| N712 | L2 | 1.0 |
| N713 | L2 | 1.0 |

Fig. 4: Seeds for the input graph

| N648 | L10 | 1.0 |
|------|------|-----|
| N649 | L10 | 1.0 |
| N650 | L10 | 1.0 |
| N651 | L10 | 1.0 |
| N652 | L10 | 1.0 |
| N653 | L10 | 1.0 |
| N654 | L10 | 1.0 |
| N655 | L10 | 1.0 |
| N656 | L10 | 1.0 |
| N657 | L10 | 1.0 |
| N658 | L10 | 1.0 |
| N659 | L10 | 1.0 |
| N660 | L10 | 1.0 |
| N661 | L10 | 1.0 |
| N662 | L10 | 1.0 |
| N663 | L10 | 1.0 |
| N664 | L10 | 1.0 |
| N665 | L10 | 1.0 |
| N666 | L10 | 1.0 |
| N667 | L10 | 1.0 |
| N668 | L10 | 1.0 |
| N669 | L10 | 1.0 |
| N670 | L10 | 1.0 |
| N671 | L10 | 1.0 |
| N672 | L10 | 1.0 |
| N673 | L10 | 1.0 |
| N674 | L10 | 1.0 |
| N675 | L10 | 1.0 |
| N676 | L10 | 1.0 |
| N677 | L10 | 1.0 |
| N678 | L10 | 1.0 |
| N679 | L10 | 1.0 |
| N680 | L10 | 1.0 |
| N681 | L10 | 1.0 |

Fig. 5: Gold labels for the unlabeled nodes.

Upon giving *input_graph, seeds* and *gold_labels* as input to *Junto,* MAD algorithm was implemented on the constructed graph and an output file, *label_prop_output* was generated. It shows the predicted label for each node with its corresponding score.

```
N28                    L1 0.06674615337582078 __DUMMY__ 0.05844988642834212 L4 0.03375105878382242 L9 0.033507898001908974 L3 0.031839068631621
4 L2 0.02536533216076376 L5 0.02531343704631428 L7 0.02478884073004222 L8 0.02475278835627636 L10 0.023949174760092323 L6 0.0237694563896475
false   0.0
N29                    L1 0.0664134634445295 __DUMMY__ 0.05999180535482096 L4 0.033403035813769036 L9 0.03288801313473841 L3 0.0313854000948446
 L2 0.02565267017002508 L7 0.025275415774037493 L8 0.02504384728583005 L10 0.024379473455828738 L5 0.02429249508684211 L6 0.02410652396224515
false   0.0
N30                    L1 0.07056983622548771 __DUMMY__ 0.05807632838515244 L4 0.03437314245428301 L9 0.03288173487773453 L3 0.0316364546032726
5 L5 0.02610290343437702 L2 0.02498104819234747 L7 0.024604270583399684 L8 0.02441588151062771 L6 0.0237586352776337 L10 0.023184638079650113
false   0.0
N31                    L1 0.06824554923149495 __DUMMY__ 0.058119657169978356 L4 0.034461922886893745 L9 0.03306410729860619 L3 0.03178368985090
588 L5 0.025931522739019185 L2 0.025251597579305018 L7 0.024886293921582783 L8 0.024235324862169377 L6 0.023769562363828833 L10 0.02322679942465
139     false   0.0
N32                    L1 0.06718820927195987 __DUMMY__ 0.05842396310834039 L4 0.034148013725134554 L9 0.033235422278360044 L3 0.03181659334765
442 L5 0.02547557180005764 L2 0.025284849323001977 L7 0.024814569712616814 L8 0.02457476745549652 L6 0.023766795338554844 L10 0.0236440522722054
7       false   0.0
N33                    L1 0.06648381677017014 __DUMMY__ 0.05806534708080696 L9 0.03374506338947878 L4 0.03300219756989735 L3 0.0329243569328051
06 L2 0.02583518685731128 L7 0.025344014225668746 L8 0.02481691515204363 L6 0.024399756041980428 L5 0.023850673187430036 L10 0.02356492101326436
8       false   0.0
N34                    L1 0.06713512178894518 __DUMMY__ 0.058116532110952676 L4 0.03323183990059806 L3 0.033176526986927364 L9 0.03252968832975
1054 L2 0.026398434254948956 L7 0.02534394198883813 L8 0.02508727308751504 L6 0.024645666372566655 L5 0.02367409290524466 L10 0.0234425909169306
2       false   0.0
N35                    L1 0.06481087984828895 __DUMMY__ 0.058766930469347216 L4 0.03372786992899605 L3 0.03293329197182736 L9 0.03249592563092
84 L2 0.025864624941455427 L7 0.025367636569746135 L8 0.0252354531147654 L6 0.024429766441194246 L5 0.02424627938377772 L10 0.0235179002526632
false   0.0
N36                    L1 0.06974604265916937 __DUMMY__ 0.05792650684737334 L4 0.033191743968548286 L9 0.03300883636941436 L3 0.032215893537293
61 L2 0.026024761122247067 L7 0.024933861660580595 L8 0.024889058864367524 L6 0.02439582458119935 L5 0.024147810111190653 L10 0.0235771393033662
4       false   0.0
N37                    L1 0.06974613741998274 __DUMMY__ 0.05787775764608731 L9 0.0341322256884788 L4 0.03226100924947576 L3 0.03195002160258446
 L2 0.025684162272090322 L8 0.025458909630194336 L7 0.025068455992710462 L5 0.024275079242320398 L6 0.023853994027464097 L10 0.0237549394174427
        false   0.0
```

Fig 6: Output file, *label_prop_output*

# RESULTS

Among the 682 predictions for the unlabeled nodes, we calculate the False Negatives (wrong predictions) and True Positives (correct predictions). Based on that the accuracy was calculated as:

True Positives = 592
False Negatives = 90

Hence, Prediction Accuracy = (True Positives) / (True Positives + False Negatives)

$$= 86.8035190616\%$$

# CONCLUSION

In our work, Modified Adsorption was implemented on a set of 982 source word – derived word pairs, and the accuracy was checked based on the no. of correct predictions. The labeled nodes (seeds) in the input graph were 300 (30% of the total no. of nodes). The trained model was 86.8% accurate, which is satisfactory.

A lot of future improvements can be done. The percentage of seed labels being input to the graph can be varied and the accuracy can be compared. The variation of accuracy can also be seen by varying th hyperparameters of the MAD implementation in *Junto*. A detailed comparison with respect to the existing label propagation algorithms can be made, and statistical measures can be derived to assess the better among the two algorithms.