

1 README:

Actor Identity through Semantic Relation Extraction

In this project a distantly supervised algorithm is implemented to gather training data that is used to train classifiers that categorize extractions of semantic relations.

1.1 Top Level Directory

Submitted with the program are the following directories used to run the program:

phase_1.cmd	This condor file calls Phase 1 of the program which encompasses several other scripts. This phase pre-processes a list of names (e.g., <code>pre-process_names.sh</code>), uses this list of names to access the html of wikipedia entries and extracts testing and training vectors (<code>create_training_data.sh</code> and <code>create_testing_data.sh</code>), creates language models(<code>create_language_model.py</code>) and extracts semantic relations (<code>extract_semantic_relations.py</code>) .
phase_2.cmd	<code>Phase_2.sh</code> is a wrapper for <code>run_semantic_extractor.py</code> pipeline. A fixed list of test sentence is sent through a pipeline to extract semantic relations without any labels and currently when called, <code>phase_2.sh</code> will target the unlabeled <code>binary_test_sentences.txt</code> file. However, any file of line separated sentences can be used. Semantic relations are extracted using the <code>get_relation()</code> function in <code>extract_semantic_relations.py</code> file
eval.cmd	For evaluation, this command calls <code>eval.sh</code> which is a wrapper for <code>./run_classifier.py</code> which is a naive bayes classifier. This is fixed to several sets of annotated data in the resource directory, <code>rec</code> , and then the <code>testing_vectors</code> directory.
rec	This is the resource directory and stores data needed to run the program. That includes language models, name lists, testing vectors, training vectors in both <code>un/annotated</code> and <code>trimmed</code> format.
src	Stores all source code to run the program. This also includes additional libraries like <code>Textacy</code> directory and <code>Cachetools</code>
output	Contains all of accuracy files that contain F1, precision, and recall scores for the binary and multi-class classifiers. Also stores system out files that show classification compared to gold standard by class for all sentences. This directory also contains a <code>sentence_file_count.txt</code> that shows the amount of relations by sentence from an extraction. Finally, there is also a <code>openie_comparison.txt</code> file that shows the Recall from the test extraction compared with the <code>openie.txt</code> file

1.2 Pre-Requisites

Python3 and NLTK must be downloaded in order to use this program.

In addition, Spacy and Textacy.extract are also required. The Spacy language model has an absolute path name to my home directory and would have to be fixed if using off of patas. Textacy is included in the `./src` directory

1.3 trimmed_data.py

Trim data takes a file of testing or training vectors and outputs a file in which class counts adhere to a set quota parameter. This script is used to ensure testing and training vectors are even because sometimes wiki extraction may be uneven.

```

├── eval.cmd
├── output
│   ├── binary_acc_10.txt
│   ├── binary_acc_25.txt
│   ├── binary_acc_50.txt
│   ├── binary_sys_output_10.txt
│   ├── binary_sys_output_25.txt
│   ├── binary_sys_output_50.txt
│   ├── binary_sys_output.txt
│   ├── class_acc_10.txt
│   ├── class_acc_25.txt
│   ├── class_acc_50.txt
│   ├── class_sys_output_10.txt
│   ├── class_sys_output_25.txt
│   ├── class_sys_output_50.txt
│   ├── class_sys_output.txt
│   ├── openie_comparison.txt
│   ├── openie.txt
│   └── sentence_file_count.txt
├── phase_1.cmd
├── phase_2.cmd
├── rec
│   ├── language_models
│   │   ├── binary_language_model.txt
│   │   └── class_language_model.txt
│   ├── name_lists
│   │   ├── name_basics.tsv
│   │   ├── testing_names.txt
│   │   └── training_names.txt
│   ├── testing_vecs
│   │   ├── binary_testing_vecs
│   │   │   ├── annotated_binary_10.txt
│   │   │   ├── annotated_binary_25.txt
│   │   │   ├── annotated_binary_50.txt
│   │   │   ├── annotated_binary.txt
│   │   │   ├── binary_testing.vectors.txt
│   │   │   ├── binary_test_relations.txt
│   │   │   └── trimmed_binary_testing.vectors.txt
│   │   ├── binary_test_sentences.txt
│   │   ├── class_testing_vecs
│   │   │   ├── annotated_class_10.txt
│   │   │   ├── annotated_class_25.txt
│   │   │   ├── annotated_class_50.txt
│   │   │   ├── annotated_class.txt
│   │   │   ├── class_testing.vectors.txt
│   │   │   ├── class_test_relations.txt
│   │   │   └── trimmed_class_testing.vectors.txt
│   │   ├── class_test_sentences.txt
│   │   └── test_semantic_relations.txt
│   └── training_vecs
│       ├── binary_training.vectors.txt
│       ├── class_training.vectors.txt
│       ├── trimmed_binary_training.vectors.txt
│       └── trimmed_class_training.vectors.txt
└── src
    ├── compare_relation_counts.py
    ├── create_language_model.py
    ├── create_testing_data.sh
    ├── create_training_data.sh
    ├── eval.sh
    ├── extract_semantic_relations.py
    ├── extract_wiki_data.py
    ├── parse_names.py
    ├── phase_1.sh
    ├── phase_2.sh
    ├── pre-process_names.sh
    ├── run_classifier.py
    ├── run_semantic_extraction.py
    └── trim_data.py

```