



**Institute for Advanced Studies
in Basic Sciences
Gava Zang, Zanjan, Iran**

NATURAL LANGUAGE PROCESSING

ONTOLOGY REPRESENTATION

Author

Sina Heydari

Master's Student of AI
at
IASBS
Institute for Advanced Studies in Basic Sciences

Spring 2024

Table Of Contents

1 A Brief Introduction To Natural Language Ontology	3
2 Ontology In Natural Language Processing	4
3 Code	5
3.1 Code Documentation	5
3.1.1 Importing The Libraries	5
3.1.2 Main Class And The Constructor	7
3.1.3 The Preprocess Method	10
3.1.4 The Build Graph Method	12
3.1.5 The Visualize Graph Method	16
3.1.6 Object Creation And Execution	19
4 Results	20
4.1 01_recurrent neural networks	20
4.2 02_recurrent neural networks	21
4.3 03_recurrent neural networks	22

4.4	04_deep learning	23
4.5	05_deep learning	24
4.6	06_deep learning	25
4.7	07_image processing	26
4.8	08_image processing	27
4.9	09_image processing	28
4.10	10_HMM	29
4.11	11_HMM	30
4.12	12_HMM	31
5	References	32

1 A Brief Introduction To Natural Language Ontology

Natural language ontology involves studying the ontological categories, structures, and concepts that are inherent in natural language. It is a branch of "descriptive metaphysics," as coined by Strawson (1959), distinct from what Fine (2017) describes as "foundational metaphysics," which focuses on what fundamentally exists. This field intersects philosophy and linguistics, specifically within metaphysics and natural language semantics. Although recognized as a distinct area only recently due to advancements in natural language semantics, philosophers have historically engaged in natural language ontology by using language to support metaphysical arguments. It differs from the ontology accepted by individuals through philosophical or naïve reasoning and from the ontology reflected in general cognition. Essentially, the ontology of a natural language represents the implicit ontological commitments made by competent speakers through their use of the language.

2 Ontology In Natural Language Processing

Understanding a domain has long been acknowledged as essential for comprehending and processing information about that domain. Ontologies, which are explicit specifications of conceptualizations, are now recognized as crucial components of information systems and information processing. This homework describes a project where ontologies are integral to the reasoning process for managing and presenting information. The processes of accessing and presenting information are facilitated through natural language, with the ontologies being integrated with the lexicon used in the natural language component.

3 Code

3.1 Code Documentation

3.1.1 Importing The Libraries

At the beginning, we import a variety of libraries and modules related to natural language processing. These imports set up the environment for a program involving processing text data, cleaning it, analyzing it, and lastly representing the results in a graph structure.

```
1 import spacy
```

Spacy: spacy is a popular library in Python for Natural Language Processing (NLP). It provides tools for processing and analyzing large volumes of text.

```
1 import glob
```

Globe: globe finds all the pathnames matching a specified pattern according to the rules used by the Unix shell. It is used for file and directory handling.

```
1 from nltk.corpus import wordnet as wn
```

WordNet: the WordNet corpus from the Natural Language Toolkit (NLTK) library aliased wn; WordNet is a lexical database for the English language, which groups words into sets of synonyms and provides short definitions and usage examples.

```
1 from spacy_cleaner import Cleaner
2 from spacy_cleaner.processing import mutators,
    removers
```

Spacy Cleaner: the spacy_cleaner is a third-party module for cleaning text data processed with spaCy. This part imports mutators and removers from the spacy_cleaner, processing submodules. These are functions or classes used to modify or remove certain parts of the text.

```
1 import networkx as nx
```

Networkx: The networkx library is imported here, which is used for creating, manipulating, and studying the structure, dynamics, and functions of graphs.

```
1 import matplotlib.pyplot as plt  
2 from matplotlib.lines import Line2D
```

Matplotlib: Generally, matplotlib is a library for visualization purposes. Here the Line2D class from matplotlib.lines is imported. Line2D is a class used for creating and manipulating lines in a plot.

3.1.2 Main Class And The Constructor

```
1 class Ontology:
```

The Ontology: class is designed to, load and process a text file stored in the given path, initializing various data structures to store and handle text data and entities, setting up a directed graph structure for potential analysis. cleaning the text using specified cleaning functions and storing the cleaned content in a dictionary. This class sets up the foundation for further text processing, analysis, and graph-based operations on the text data. A class in Python is a blueprint for creating objects (a particular data structure), providing initial values for state (member variables), and implementations of behavior (member functions or methods).

```
1     def __init__(self, path):
```

__init__: This is the constructor method for the Ontology class. It initializes an instance of the class. The self parameter refers to the instance being created, and path is an argument that specifies the file path to a the text file the code want to work on.

```
1         self.path = path
2         self.language_model =
3             spacy.load("en_core_web_sm")
```

- **self.path:** stores the file path provided during the creation of the Ontology object.
- **self.language_model:** loads an English language model from spaCy, used for NLP tasks.

```

1 # self.graph = nx.MultiDiGraph()
2 self.graph = nx.DiGraph()
3 # self.graph = nx.Graph()

```

self.graph: initializes a directed graph using NetworkX's DiGraph class, which allows us to do the representation with directed graphs. The commented lines suggest alternatives (MultiDiGraph for a directed graph that can hold multiple edges between nodes, and Graph for an undirected graph).

```

1 self.corpus_uniques = {}
2 self.corpus_dict = {}
3 self.entity_list = []
4 self.corpus_list = []

```

These lines initialize several empty data structures:

- **self.corpus_uniques:** A dictionary intended to store unique elements from the corpus.
- **self.corpus_dict:** A dictionary to store the processed text content.
- **self.entity_list:** A list to store entities extracted from the text.
- **self.corpus_list:** A list to store the corpus or multiple corpora.

```

1 self.lemma = Cleaner(self.language_model,
2                     removers.remove_number_token,
3                     removers.remove_punctuation_token,
4                     mutators.mutate_lemma_token)

```

self.lemma: initializes a Cleaner object from the spacy_cleaner module, configured with the spaCy language model and several cleaning functions:

- **removers.remove_number_token:** Function to remove number tokens.
- **removers.remove_punctuation_token:** Function to remove punctuation tokens.
- **mutators.mutate_lemma_token:** Function to change tokens to their lemmas(lemmatizing the tokens).

```

1      with open(path, 'r', encoding="latin-1") as
2          file:
3              content = file.read()
4              # Remove unwanted characters and
5              # normalize whitespace
6              content = content.replace('\r', '') \
7                  .replace('\n', ' ') \
8                  .replace('\x92', '') \
9                  .replace('\x93', '') \
10                 .replace('\x94', '') \
11                 .replace('\x95', '') \
12                 .replace('\x96', '') \
13                 .replace('\x97', '') \
14                 .replace('\x98', '') \
15                 .replace('?', ' ') \
16                 .replace('!', ' ') \
17                 .replace('.', ' ') \
18                 .replace(':', ' ') \
19                 .replace(';', ' ')

```

This block opens the file specified by path for reading with the encoding latin-1. It reads the file content into the variable content. The replace method is used to clean the text by removing unwanted characters and normalizing whitespace. This includes removing carriage returns, newline characters, and various non-standard punctuation characters, and replacing punctuation marks with spaces.

```

1      self.corpus_dict[0] = content

```

Finally, the clean text is stored in the self.corpus_dict dictionary with the key 0.

3.1.3 The Preprocess Method

```
1 def preprocess(self):
```

Preprocess: This line defines a **method** called preprocess within the Ontology class. This method will preprocess the text stored in self.corpus_dict. The preprocess method performs several steps to prepare the text data for analysis. Including, all text conversion to lowercase to ensure uniformity, text cleaning routines to remove unwanted elements and standardize tokens, stores the cleaned text in a list for easy processing, annotating the text with linguistic features using the spaCy language model, extracting unique, non-stopword tokens and mapping them to their WordNet synsets, and building a dictionary of unique tokens as keys, and their synsets as values. This method effectively prepares the text for further analysis, such as entity recognition, relationship extraction, and graph-based modeling.

```
1 # Convert text to lowercase
2 for key in self.corpus_dict:
3     self.corpus_dict[key] =
        self.corpus_dict[key].lower()
```

This **loop** iterates over each key in self.corpus_dict. For each key, it converts the corresponding text to lowercase to standardize the text and make the subsequent processing case-insensitive.

```
1 # Clean text using spacy_cleaner
2 for key in self.corpus_dict:
3     self.corpus_dict[key] =
        self.lemma.clean([str(self.corpus_dict[key])])
```

This **loop** again iterates over each key in self.corpus_dict. For each key, it applies the self.lemma.clean method to clean the text using the Cleaner object initialized earlier. This involves removing numbers and punctuation and converting tokens to their lemmas

```

1   # Store cleaned text in a list
2   for key in self.corpus_dict:
3       self.corpus_list.append(self.corpus_dict[key])

```

This **loop** iterates over each key in `self.corpus_dict`. For each key, it appends the cleaned text to `self.corpus_list`, creating a list of cleaned text entries

```

1   # Process text with spaCy language model
2   for i in range(len(self.corpus_list)):
3       self.corpus_list[i] =
            self.language_model(str(self.corpus_list[i]))

```

This **loop** iterates over each index in `self.corpus_list`. For each index, it processes the text using the spaCy language model (`self.language_model`). This converts each text entry into a spaCy Doc object, which contains linguistic annotations.

```

1   # Extract unique non-stop tokens and find their
2   # WordNet synsets
3   for i in range(len(self.corpus_list)):
4       tokens = self.corpus_list[i]
5       for token in tokens:
6           if not token.is_stop:
7               if str(token) not in
                    self.corpus_uniques:
                        self.corpus_uniques[str(token)] =
                            set(wn.synsets(str(token)))

```

This **nested loop** iterates over each Doc object in `self.corpus_list`. For each Doc object, it iterates over each token in the Doc. For each token, it checks if the token is not a stop word (`token.is_stop`). If the token is not a stop word and is not already in `self.corpus_uniques`, it adds the token to `self.corpus_uniques` with a set of its corresponding WordNet synsets. This maps each unique token to its potential meanings as defined in WordNet.

3.1.4 The Build Graph Method

```
1 def build_graph(self):
```

The **build_graph** is a **method** within the **Ontology class**. This method will construct a graph based on the unique tokens and their relationships found in WordNet. The **build_graph method** constructs a graph by:

- **Nodes:** Each unique token from the preprocessed text is added as a node.
- **Edges:** Edges are added between nodes based on semantic relationships found in WordNet, including synset overlap, hypernym, hyponym, synonym, antonym, meronym, and holonym relationships.
- **Relationship Detection:** The method checks for each type of relationship and creates edges with appropriate attributes, building a semantically rich graph of the text data.

The graph, will further be used for ontology representation of the given documents.

```
1 # Add unique tokens as nodes in the graph
2 for node in self.corpus_uniques.keys():
3     self.graph.add_node(node)
```

This loop iterates over each unique token in `self.corpus_uniques`. For each token, it adds the token as a node to the graph `self.graph`.

```
1 nodes = list(self.corpus_uniques.keys())
2 # Add edges between nodes based on
3     # relationships found in WordNet
4 for i in range(len(nodes)):
5     for j in range(i + 1, len(nodes)):
6         node1 = nodes[i]
7         node2 = nodes[j]
8         synsets1 = self.corpus_uniques[node1]
9         synsets2 = self.corpus_uniques[node2]
10
11         relationship_found = False
12         edge_attributes = {}
```

nodes is a list of all unique tokens (graph nodes). Two nested loops iterate over pairs of nodes. "i" is the index of the first node and "j" is the index of the second node. For each pair of nodes (node1 and node2), it retrieves their corresponding synsets (synsets1 and synsets2). relationship_found is a flag to track if a relationship is found between the nodes. And edge_attributes is a dictionary to store attributes of the edge if a relationship is found.

```

1      # Check for direct synset overlap
2      if synsets1 & synsets2:
3          relationship_found = True
4          edge_attributes['type'] =
5              'synset_overlap'
```

This part checks if there is a direct overlap between synsets1 and synsets2. If an overlap is found, it sets **relationship_found** to **True** and adds an attribute '**type**':

- **synset_overlap** to **edge_attributes**.

```

1      else:
2          # Check for other meaningful
3          # relationships
4          for syn1 in synsets1:
5              for syn2 in synsets2:
6                  if syn1 in syn2.hypernyms()
7                      or syn2 in
8                          syn1.hypernyms():
9                          relationship_found =
10                             True
11                         edge_attributes['type']
12                         = 'hypernym'
13                         elif syn1 in
14                             syn2.hyponyms() or syn2
15                             in syn1.hyponyms():
16                             relationship_found =
17                                 True
18                                 edge_attributes['type']
19                                 = 'hyponym'
20                                 elif any(lemma in
21                                     syn2.lemmas() for lemma
22                                     in syn1.lemmas()):
```

```

12         relationship_found =
13             True
14             edge_attributes[ 'type' ]
15                 = 'synonym'
16             elif any(antonym for lemma
17                 in syn1.lemmas() \
18                     for antonym in
19                         lemma.antonyms() \
20                             if antonym in
21                                 syn2.lemmas()): :
22             relationship_found =
23                 True
24                 edge_attributes[ 'type' ]
25                     = 'antonym'
26             elif syn1 in
27                 syn2.part_meronyms() or
28                     syn2 in
29                         syn1.part_meronyms():
30                         relationship_found =
31                             True
32                             edge_attributes[ 'type' ]
33                                 = 'meronym'
34             elif syn1 in
35                 syn2.part_holonyms() or
36                     syn2 in
37                         syn1.part_holonyms():
38                         relationship_found =
39                             True
40                             edge_attributes[ 'type' ]
41                                 = 'holonym'
42
43             if relationship_found:
44                 self.graph.add_edge(node1,
45                     node2,
46                         **edge_attributes)
47                     break
48             if relationship_found:
49                 break

```

This **nested loop** iterates over each synset in synsets1 and synsets2. For each pair of synsets (syn1 and syn2), it checks for various types of relationships:

- **Hypernym:** If syn1 is a hypernym of syn2 or vice versa.
- **Hyponym:** If syn1 is a hyponym of syn2 or vice versa.
- **Synonym:** If any lemma in syn1 matches a lemma in syn2.
- **Antonym:** If any lemma in syn1 has an antonym in syn2.
- **Meronym:** If syn1 is a part meronym of syn2 or vice versa.
- **Holonym:** If syn1 is a part holonym of syn2 or vice versa.

If any relationship is found, it sets **relationship_found** to **True**, assigns the appropriate relationship type to **edge_attributes**, and adds an **edge** between **node1** and **node2** with the attributes.

The break statements exit the inner loops once a relationship is found to avoid redundant checks.

3.1.5 The Visualize Graph Method

```
1 def visualize_graph(self):
```

This method provides a visual representation of the semantic relationships between the unique tokens extracted from the document, allowing for easier interpretation and analysis.

The **method** visualizes the graph of entities and their relationships; cleans the graph by removing nodes without edges; uses a spring layout for positioning nodes; sets edge colors based on the types of relationships; plots the graph with labels, node sizes, colors, and edge colors; creates and displays a legend for the edge colors. sets the title and displays the graph.

```
1     # Remove nodes without any edges
2     isolated_nodes = list(nx.isolates(self.graph))
3     self.graph.remove_nodes_from(isolated_nodes)
```

This line defines a **method visualize_graph** within the Ontology class. This method will visualize the graph constructed in the build_graph method.

```
1     # Define layout for graph visualization
2     pos = nx.spring_layout(self.graph, k=1,
3                             iterations=50)
4     plt.figure(figsize=(15, 15))
```

This line defines the layout for the graph visualization using the `spring_layout` function from **NetworkX**. The `spring_layout` positions nodes using the Fruchterman-Reingold force-directed algorithm.

- **k=1** is the optimal distance between nodes, and `iterations=50` is the number of iterations to perform.
- `plt.figure(figsize = (15, 15))` sets the figure size for the plot to 15x15 inches.

```

1          # Define edge colors based on
2          # relationship types
3
4      edge_colors = []
5      edge_color_map = {
6          'synset_overlap': 'black',
7          'hypernym': 'blue',
8          'hyponym': 'green',
9          'synonym': 'purple',
10         'antonym': 'red',
11         'meronym': 'orange',
12         'holonym': 'brown'
13     }
14
15     for _, _, data in self.graph.edges(data=True):
16         edge_colors.append(edge_color_map.get(data['type'],
17                                         'gray'))

```

This section initializes an empty list `edge_colors` to store the colors of edges. `edge_color_map` is a dictionary that maps relationship types to specific colors. The loop iterates over all edges in the graph, retrieving the data associated with each edge. For each edge, it appends the appropriate color (based on the relationship type) to `edge_colors`. If the relationship type is not found in `edge_color_map`, it defaults to `gray`.

```

1          # Draw the graph with labels,
2          # colors, and sizes
3
4      nx.draw(self.graph, pos, with_labels=True,
5              node_size=2500,
6              node_color="yellow", font_size=10,
7              font_weight="bold",
8              edge_color=edge_colors)

```

This line uses the `nx.draw` function to draw the graph with various parameters.

pos: Node positions defined by the spring layout.

- `with_labels=True`: Display node labels.
- `node_size=2500`: Set the size of the nodes.
- `node_color="yellow"`: Set the color of the nodes.

- `font_size=10`: Set the font size of the labels.
- `font_weight="bold"`: Set the font weight of the labels.
- `edge_color=edge_colors`: Set the color of the edges based on their relationship types

```

1 # Create legend for edge colors
2 legend_elements = [Line2D([0], [0],
3                         color=color, linewidth=4, label=label)
4                     for label, color in
5                         edge_color_map.items()]
6 plt.legend(handles=legend_elements, loc='best')

```

This section creates a legend for the edge colors using `Line2D` from `matplotlib.lines`. `legend_elements` is a list of `Line2D` objects, each representing a different relationship type and its corresponding color. `plt.legend(handles=legend_elements, loc='best')` adds the legend to the plot at the best location.

```

1 # Set the title and display the graph
2 plt.title(f'Document_{self.path}_Entity_Graph')
3 plt.show()

```

This line sets the title of the plot to indicate which document the graph represents, using the path of the document. `plt.show()` displays the graph.

3.1.6 Object Creation And Execution

```
1 # Representing one document
2 document_list = glob.glob('*.txt')
3
4 represent = Ontology(document_list[0])
5 represent.preprocess()
6 represent.build_graph()
7 represent.visualize_graph()
8
9 # Representing all the documents (in separate
10 # figures)
11 for document in document_list:
12     represent = Ontology(document)
13     represent.preprocess()
14     represent.build_graph()
15     represent.visualize_graph()
```

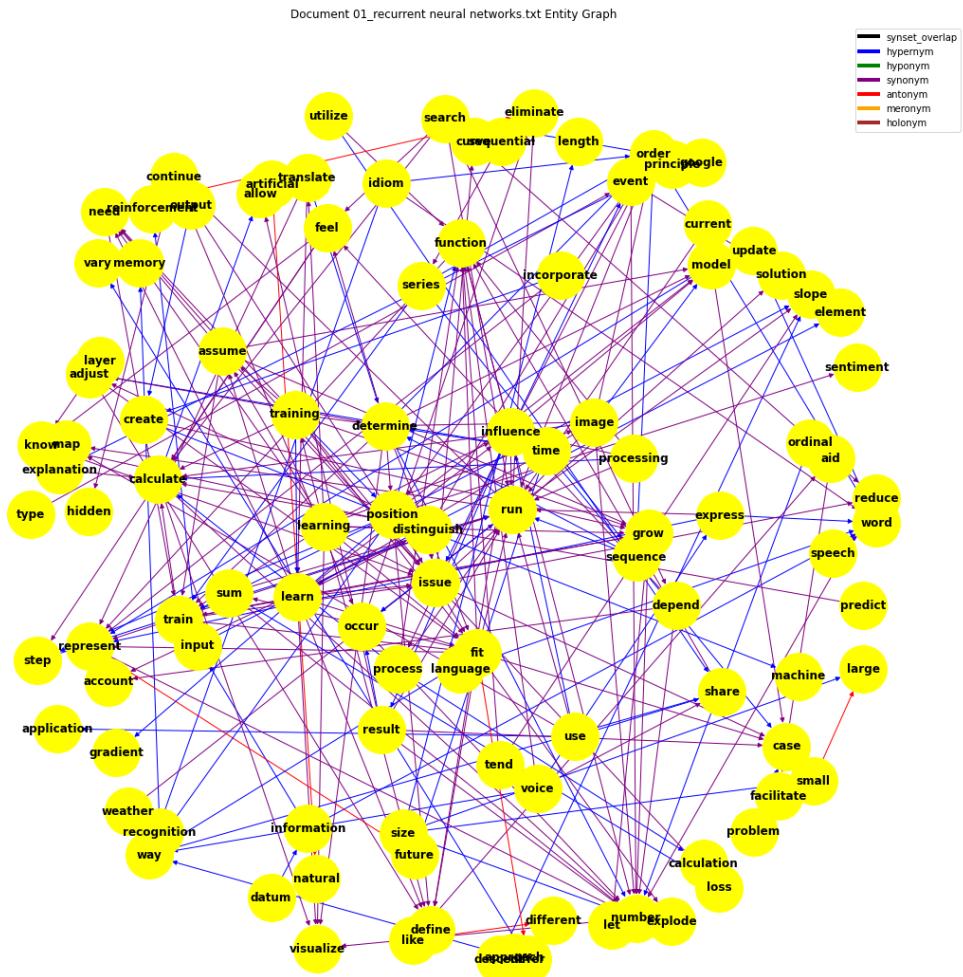
This code snippet demonstrates two main functionalities:

- **Single Document Processing and Visualization**, Processes and visualizes the graph for the first document in the list.
- **Multiple Document Processing and Visualization**, Iterates through each document in the list, processes, and visualizes the graph for each document in separate figures.

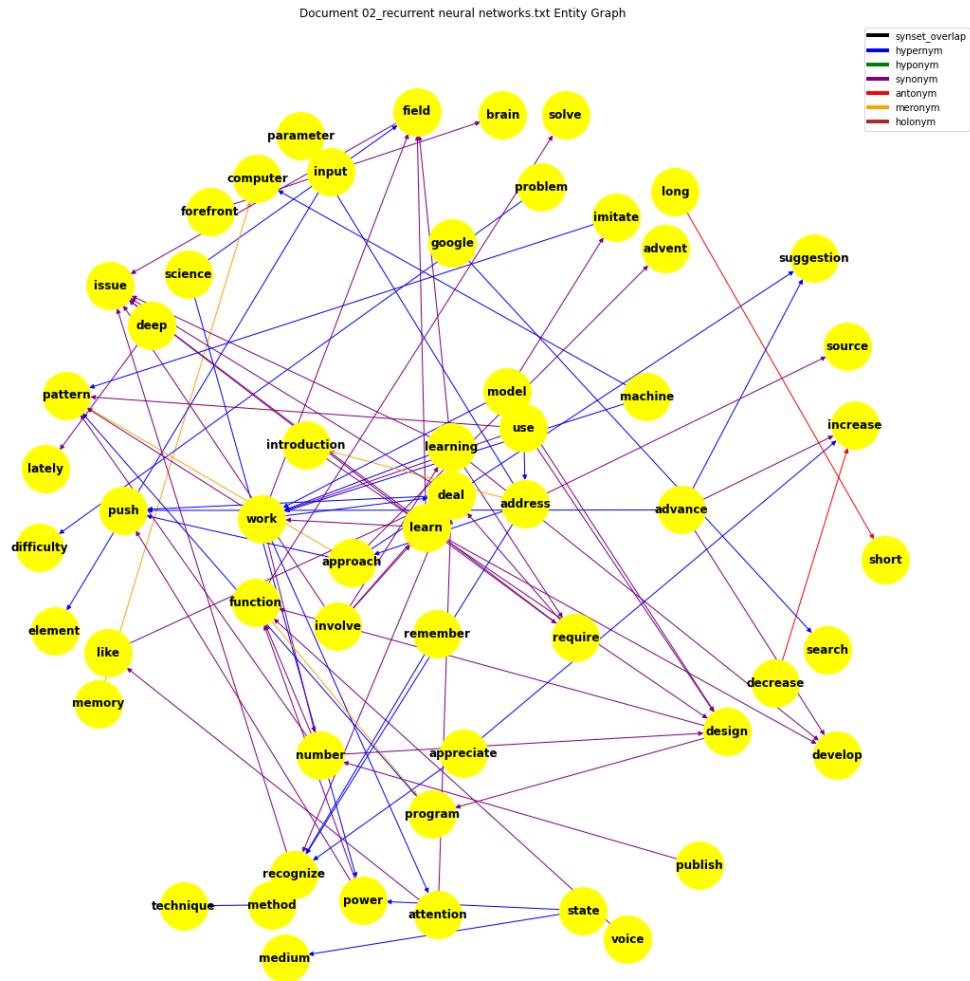
This part allows for detailed semantic analysis of individual documents and a comprehensive analysis of multiple documents by visualizing their entity relationships.

4 Results

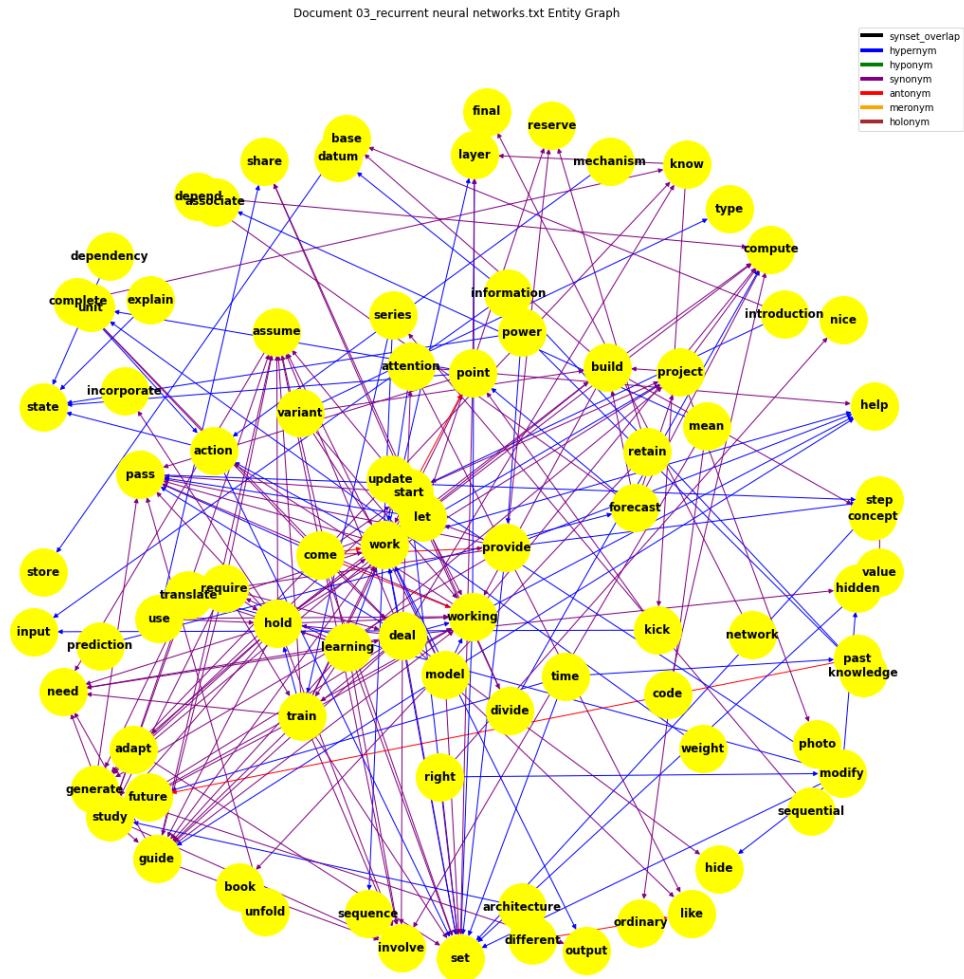
4.1 01_recurrent neural networks



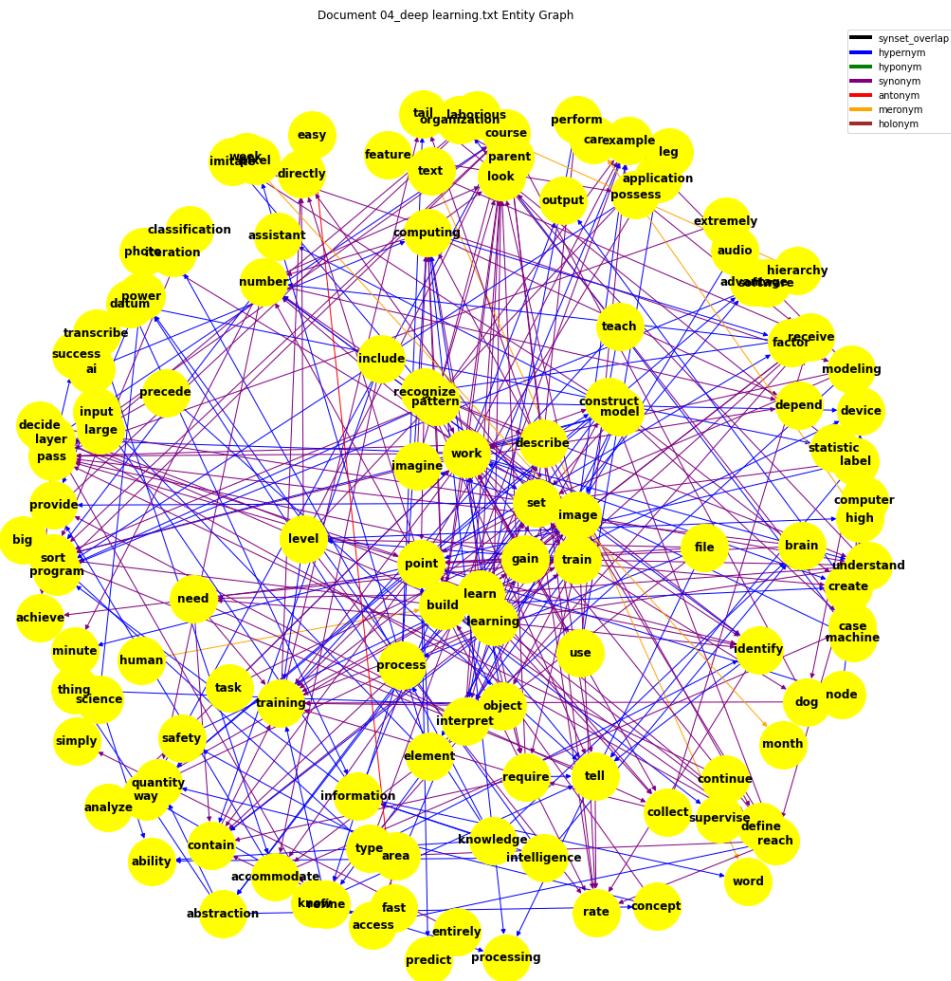
4.2 02_recurrent neural networks



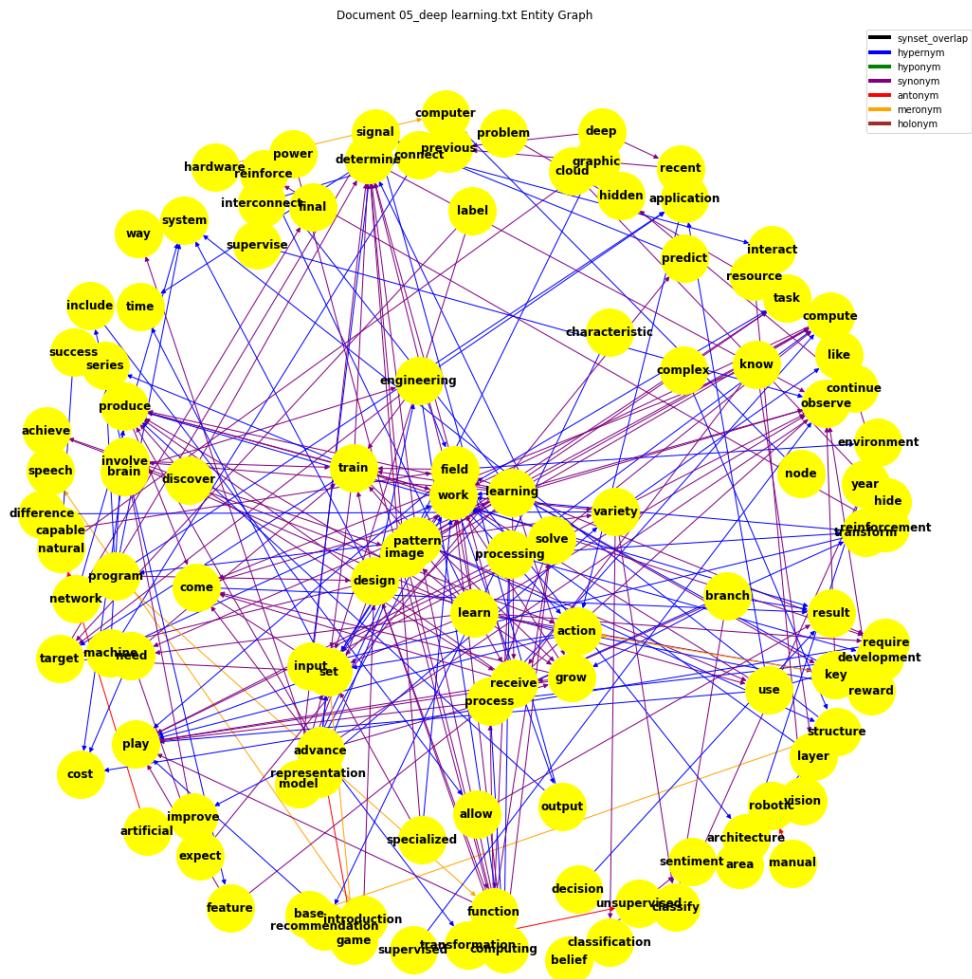
4.3 03_recurrent neural networks



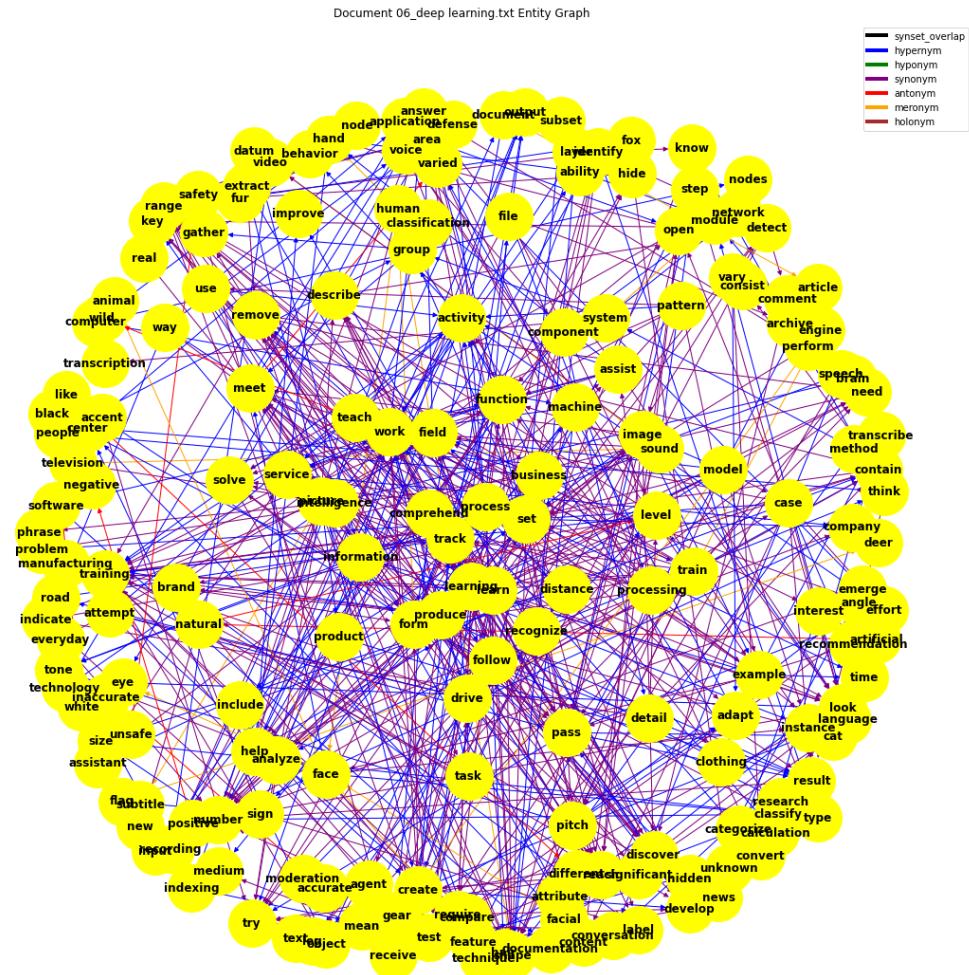
4.4 04_deep learning



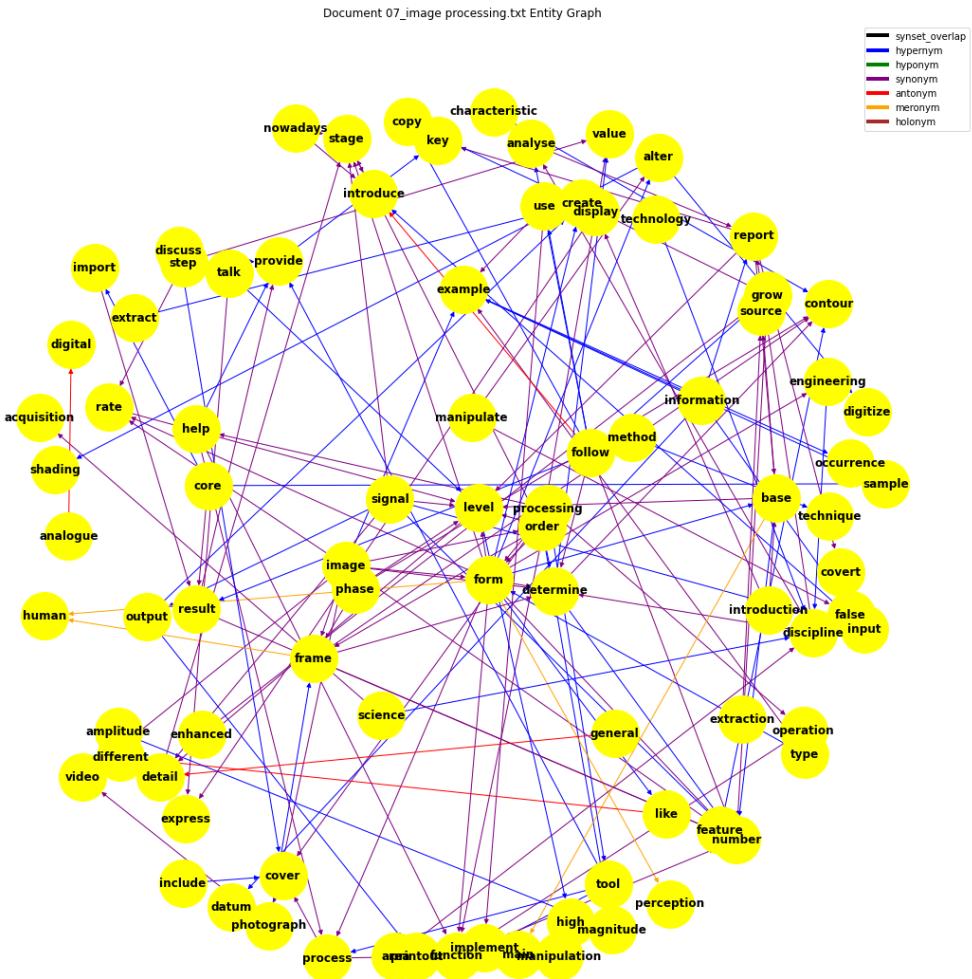
4.5 05_deep learning



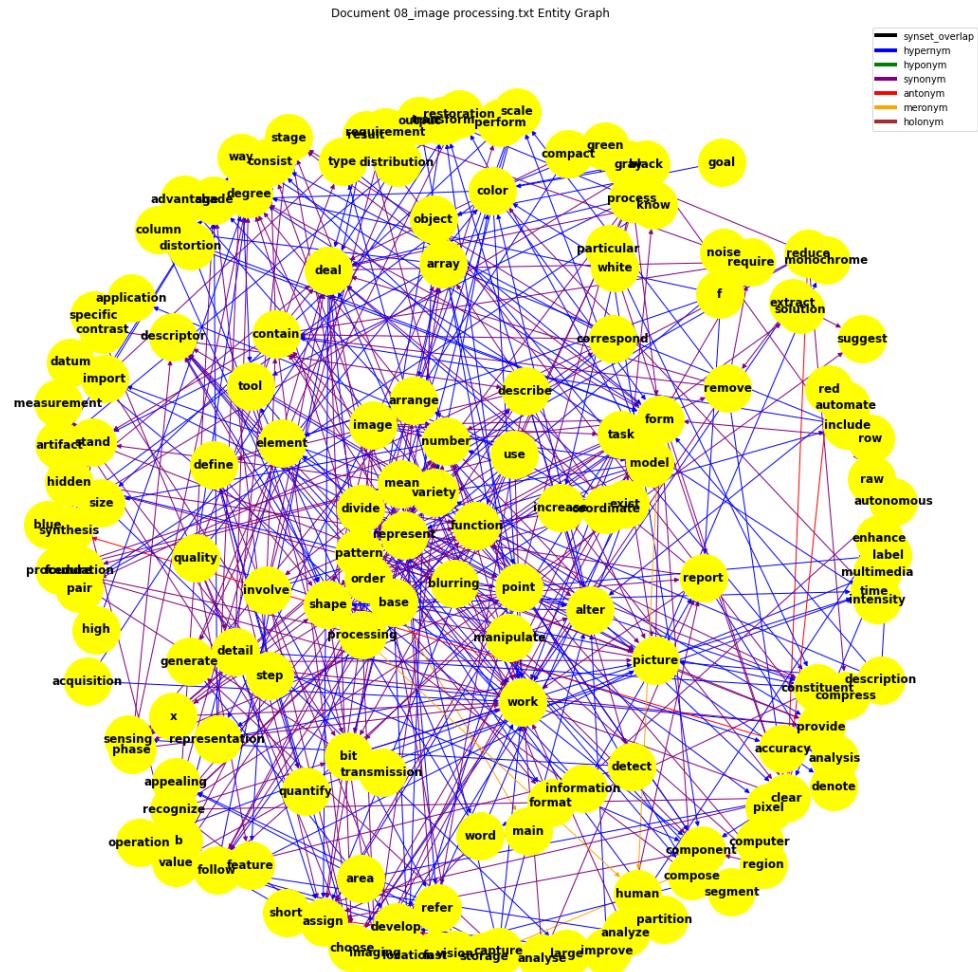
4.6 06_deep learning



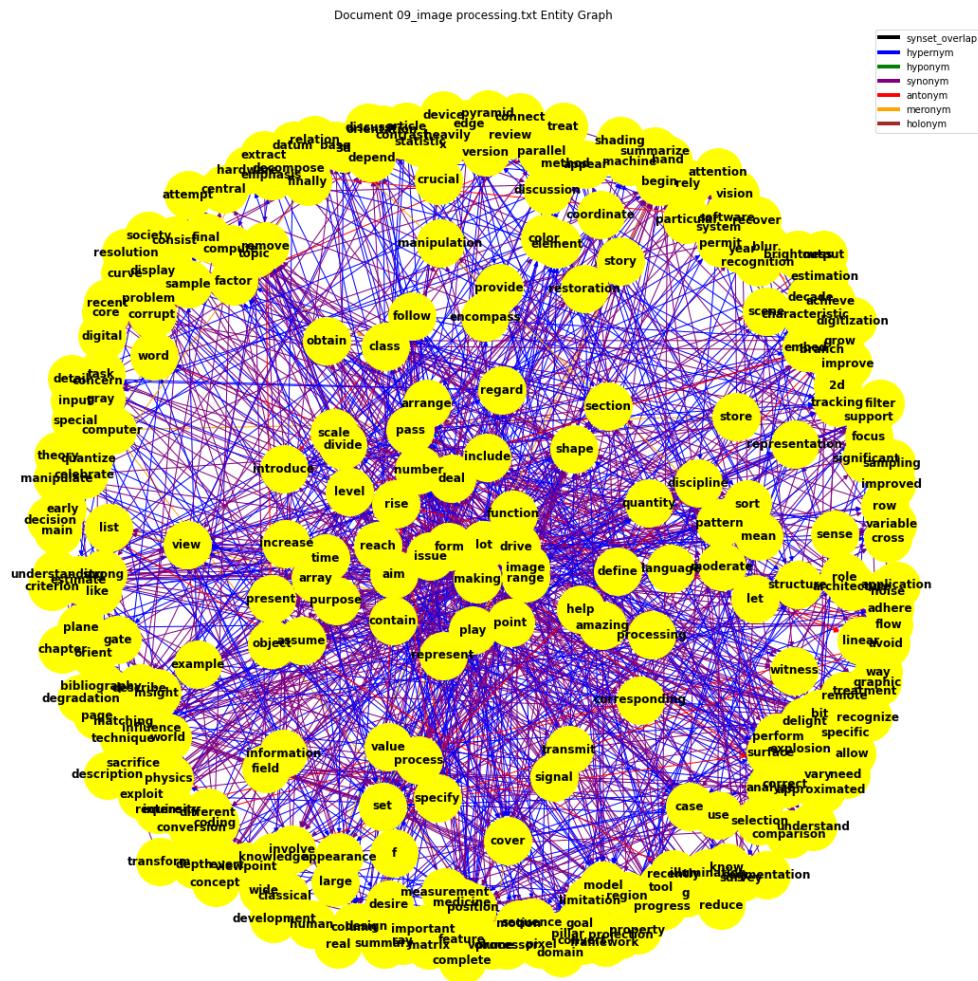
4.7 07_image processing



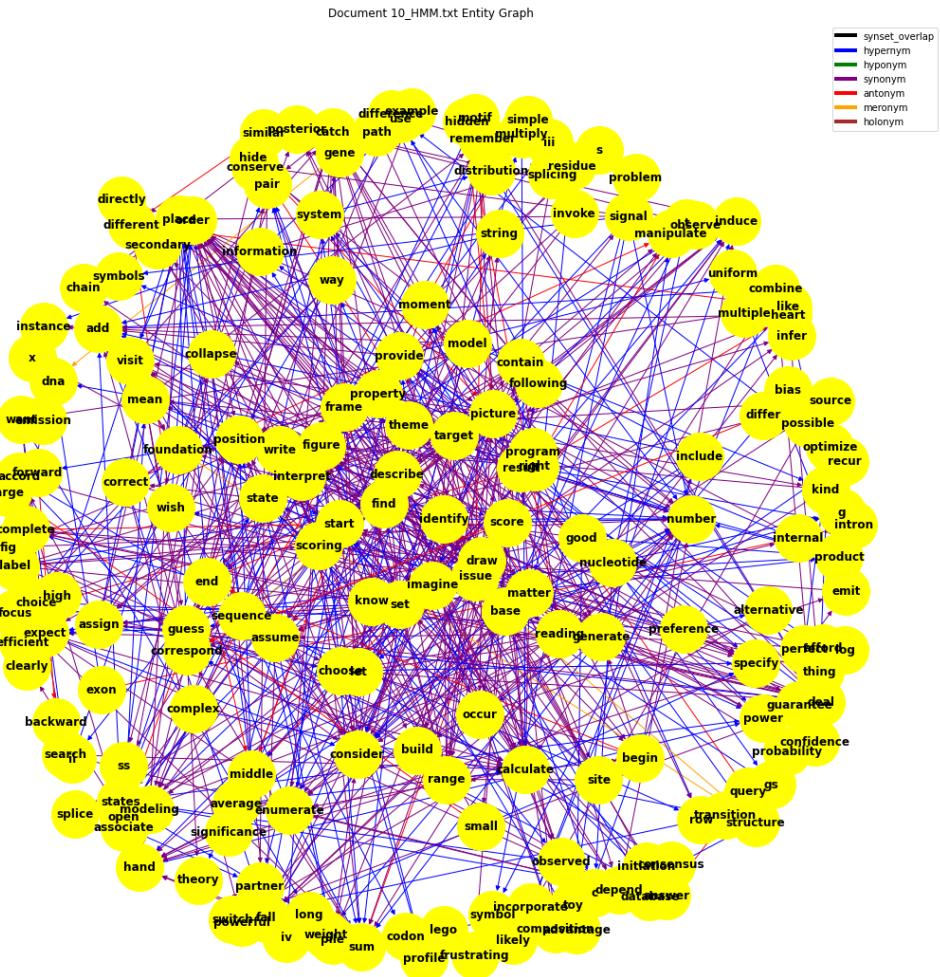
4.8 08_image processing



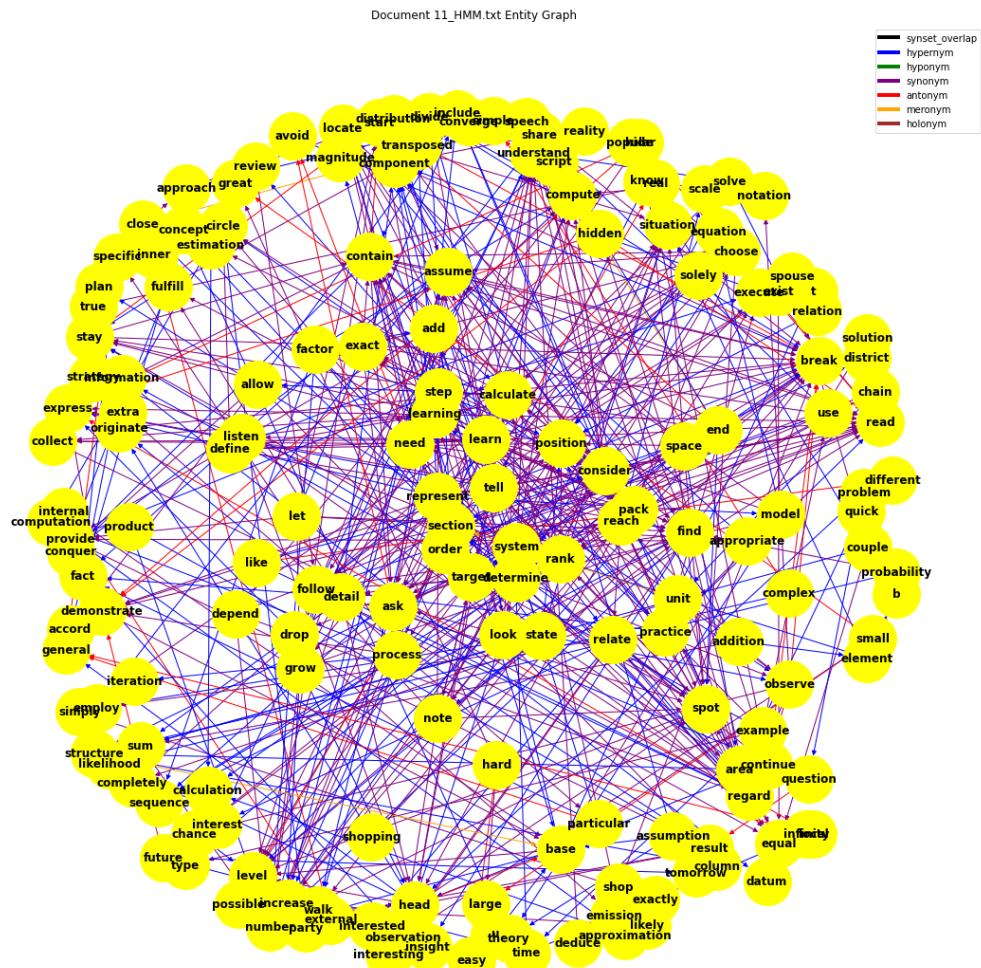
4.9 09_image processing



4.10 10_HMM

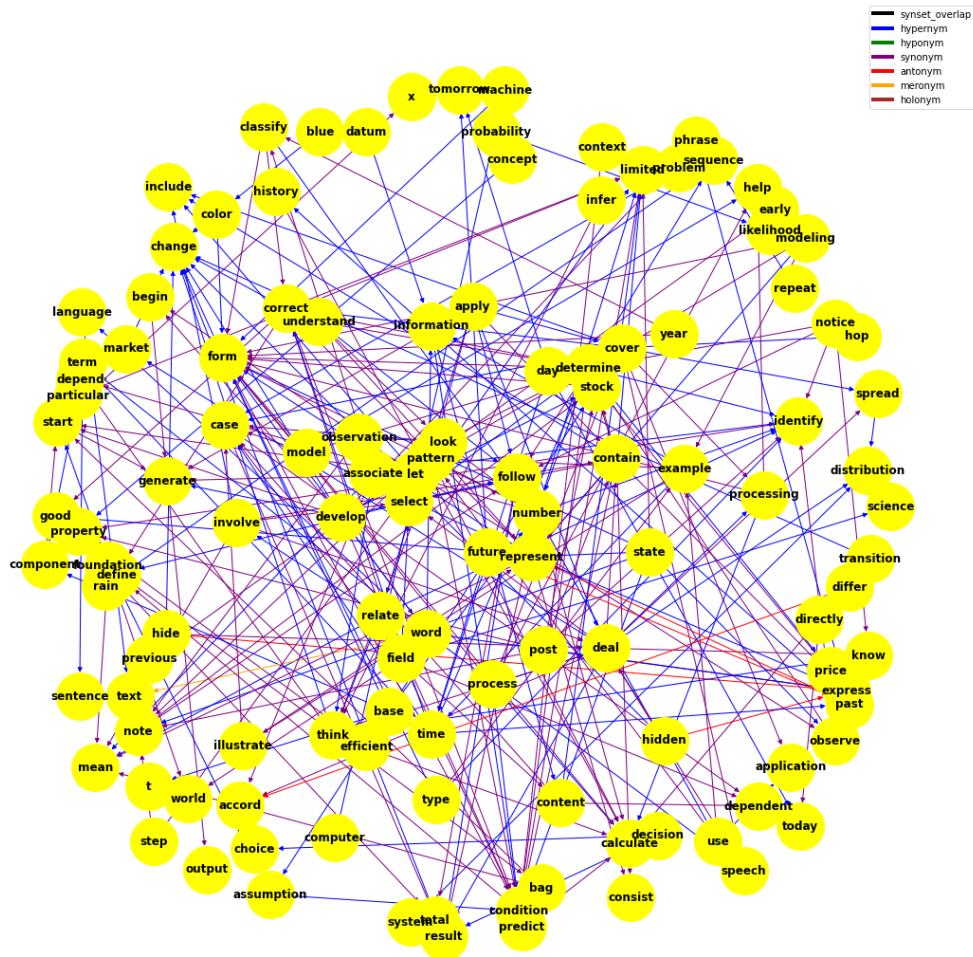


4.11 11_HMM



4.12 12_HMM

Document 12_HMM.txt Entity Graph



5 References

- Ontology In Natural Language
- Introduction to Ontologies
- Semantic Web and RDF
- spacy-cleaner
- spacy
- glob
- NetworkX
- Matplotlib
- nltk