

Lecture Information

Course: CSCI E-89B: Natural Language Processing
Lecture: Lecture 10
Topic: Named Entity Recognition (NER)
Date: Fall 2024

Contents

1 Introduction to Named Entity Recognition

Overview

Named Entity Recognition (NER) is the task of identifying and classifying named entities in text into predefined categories such as person names, organizations, locations, dates, and more. It's a fundamental NLP task with applications in information extraction, question answering, and machine translation.

1.1 What are Named Entities?

Named Entity Categories

Common named entity types include:

- **PERSON**: Names of people (Donald Trump, Marie Curie)
- **ORG**: Organizations, companies, institutions (Tesla, Harvard University)
- **GPE**: Geopolitical entities—countries, cities, states (America, Paris)
- **LOC**: Non-GPE locations (Mount Everest, Pacific Ocean)
- **DATE**: Dates and time periods (November 4, 2024, this year)
- **TIME**: Times (3:00 PM, midnight)
- **MONEY**: Monetary values (\$1 trillion, 50 euros)
- **PERCENT**: Percentages (40%, two-thirds)
- **CARDINAL**: Numbers not fitting other categories (50, three)
- **ORDINAL**: Ordinal numbers (first, 2nd)
- **NORP**: Nationalities, religious/political groups (Chinese, Republican)

1.2 Why is NER Important?

Important

NER is crucial for many downstream tasks:

- **Machine Translation**: Knowing “Tesla” is an organization helps translate correctly
- **Information Extraction**: Extract structured data from unstructured text
- **Question Answering**: Identify entities mentioned in questions
- **Search**: Improve semantic search by understanding entity types
- **Sentiment Analysis**: Attribute sentiment to specific entities

NER in Action

Input text: “Donald Trump won more than 50 electoral votes this year. Tesla’s stock rose 2.4%.”

NER output:

Entity	Type	Position
Donald Trump	PERSON	0–11
more than 50	CARDINAL	17–29
this year	DATE	47–56
Tesla	ORG	58–63
2.4%	PERCENT	78–82

1.3 NER for Feature Enhancement

Enhancing Classification with NER

NER can improve text classification by:

1. **Adding entity counts:** Concatenate counts of persons, organizations, etc.
2. **Entity-based features:** Create binary indicators for entity presence
3. **Structured metadata:** Extract entities as document metadata
4. **Relationship extraction:** Find connections between entities

2 Two Approaches to NER

Overview

NER can be performed using two main approaches: rule-based methods using pattern matching, and statistical/neural methods using machine learning.

2.1 Approach Comparison

Aspect	Rule-Based	Statistical/Neural
Training Data	Not required	Required (labeled)
Context Awareness	Limited	High
Maintenance	High burden	Lower
Adaptability	Poor	Good
Interpretability	High	Lower
Scalability	Poor	Good
Accuracy	Depends on rules	Generally higher

3 Rule-Based NER

Overview

Rule-based NER uses handcrafted patterns (regular expressions, dictionaries, linguistic rules) to identify entities. While limited, it's interpretable and requires no training data.

3.1 Regular Expressions for Entity Detection

Common Pattern Components

- \b : Word boundary
- \d{n} : Exactly n digits
- \d{1,2} : 1 or 2 digits
- \s+ : One or more whitespace characters
- [A-Z] : One uppercase letter
- [a-z]+ : One or more lowercase letters
- (?:...) : Non-capturing group
- ?: Makes preceding element optional
- : OR operator

3.2 Date Pattern Example

```

1 import re
2
3 # Pattern for dates like "11/04/2024" or "November 4, 2024"
4 date_pattern = r'',
5     \b                      # Word boundary
6    (?::
7         \d{1,2}/\d{1,2}/\d{4}      # MM/DD/YYYY format
8         |
9        (?:January|February|March|April|May|June|
10        July|August|September|October|November|December)
11        \s+                     # Required space
12        \d{1,2}                  # Day (1-31)
13       (?: ,\s*)?              # Optional comma and space
14        \d{4}                   # Year
15    )
16    \b
17 '',
18
19 text = "The event is on November 4, 2024 or 11/04/2024."
20 dates = re.findall(date_pattern, text, re.VERBOSE)
21 print(dates) # ['November 4, 2024', '11/04/2024']

```

Listing 1: Regular Expression for Dates

3.3 Person Name Pattern

```
1 # Pattern for names with titles
```

```

2 person_pattern = r''''
3     \b
4     (?:(Mr\.|Mrs\.|Ms\.|Dr\.|Professor)    # Title
5     \s+                                # Space
6     [A-Z][a-z]+                      # First name (capitalized)
7     (?:(\s+[A-Z][a-z]+)?            # Optional last name
8     \b
9 '''
10
11 text = "Mr. Trump met with Dr. Smith yesterday."
12 persons = re.findall(person_pattern, text, re.VERBOSE)
13 print(persons)  # ['Mr. Trump', 'Dr. Smith']

```

Listing 2: Pattern for Names with Titles

3.4 Organization Pattern

```

1 # Pattern for company names
2 org_pattern = r''''
3     \b
4     [A-Z][a-zA-Z\s]+                  # Company name
5     (?:(Inc\.|Ltd\.|Corporation|Corp\.)  # Corporate suffix
6     \b
7 '''
8
9 text = "Apple Inc. announced new products."
10 orgs = re.findall(org_pattern, text, re.VERBOSE)
11 print(orgs)  # ['Apple Inc.']

```

Listing 3: Pattern for Organizations

3.5 Complete Rule-Based NER System

```

1 import re
2
3 def rule_based_ner(text):
4     entities = []
5
6     # Date patterns
7     date_patterns = [
8         r'\b\d{1,2}/\d{1,2}/\d{4}\b',
9         r'\b(?:Jan|Feb|Mar|Apr|May|Jun|Jul|Aug|Sep|Oct|Nov|Dec)[a-z]*\s+\d{1,2},?\s*\d{4}\b'
10    ]
11
12     # Email pattern
13     email_pattern = r'\b[\w.]+@[^\w.]+\.\w+\b'
14
15     # Time pattern
16     time_pattern = r'\b\d{1,2}:\d{2}\s*(?:AM|PM|am|pm)?\b'
17
18     # Person pattern (with titles)
19     person_pattern = r'\b(?:Mr\.|Mrs\.|Ms\.|Dr\.)\s+[A-Z][a-z]+(?:(\s+[A-Z][a-z]+)?\b'
20
21     # Organization pattern
22     org_pattern = r'\b[A-Z][a-zA-Z\s]+(?:(Inc\.|Ltd\.|Corp\.)\b'
23
24     # Apply patterns
25     for pattern in date_patterns:
26         for match in re.finditer(pattern, text):

```

```

27     entities.append((‘DATE’, match.group(), match.span()))
28
29     for match in re.finditer(email_pattern, text):
30         entities.append((‘EMAIL’, match.group(), match.span()))
31
32     for match in re.finditer(time_pattern, text):
33         entities.append((‘TIME’, match.group(), match.span()))
34
35     for match in re.finditer(person_pattern, text):
36         entities.append((‘PERSON’, match.group(), match.span()))
37
38     for match in re.finditer(org_pattern, text):
39         entities.append((‘ORG’, match.group(), match.span()))
40
41     return entities
42
43 # Test
44 text = """
45 Meeting with Dr. Smith at 3:00 PM on November 18, 2024.
46 Contact: john.doe@company.com. Apple Inc. will attend.
47 """
48 entities = rule_based_ner(text)
49 for entity_type, entity, span in entities:
50     print(f"{entity_type}: {entity} at {span}")

```

Listing 4: Simple Rule-Based NER System

3.6 Limitations of Rule-Based NER

Rule-Based Limitations

- **No context awareness:** “Tesla” could be a person (Nikola Tesla) or company
- **Missing entities:** “Donald Trump” without title won’t be recognized
- **Language-specific:** Rules must be rewritten for each language
- **Date format variations:** US vs European formats differ
- **Maintenance burden:** Rules must be constantly updated
- **Name variations:** “La Place” vs “Laplace” requires special handling

4 Statistical and Neural NER

Overview

Modern NER systems use machine learning to learn patterns from labeled data. Neural networks, particularly CNNs and transformers, achieve state-of-the-art performance.

4.1 NER as Sequence Labeling

BIO Tagging Scheme

NER is typically framed as a sequence labeling problem using BIO tags:

- **B-TYPE:** Beginning of an entity of TYPE
- **I-TYPE:** Inside/continuation of an entity

- **O:** Outside any entity

BIO Tagging Example

Sentence: “Donald Trump visited Tesla headquarters.”

Token	Tag
Donald	B-PERSON
Trump	I-PERSON
visited	O
Tesla	B-ORG
headquarters	O
.	O

4.2 Statistical Methods

Traditional ML for NER

Hidden Markov Models (HMM):

- Model sequence of tags as Markov chain
- Assume current tag depends only on previous tag
- Limited feature representation

Conditional Random Fields (CRF):

- Discriminative model (directly models $P(\text{tags}|\text{words})$)
- Can use arbitrary features
- Global normalization (considers entire sequence)

4.3 Neural Network Architecture for NER

Important

SpaCy uses a **Convolutional Neural Network (CNN)** for NER. Here's why:

1. Word embeddings capture semantic meaning
2. CNN filters capture local context (neighboring words)
3. Sliding window naturally handles variable-length text
4. Efficient parallel computation

4.4 How Neural NER Works

Neural NER Pipeline

1. **Input:** Document text
2. **Tokenization:** Split into tokens
3. **Embedding:** Convert tokens to vectors (e.g., 4D, 100D)
4. **CNN:** Apply filters over embedding sequences
5. **Output Layer:** Softmax over entity types for each token

NER as CNN Classification

Input: “The cat sat on mat”

Step 1: Embed each token:

Token	Dim 1	Dim 2	Dim 3	Dim 4
The	0.8	0.1	0.3	0.7
cat	0.5	0.7	0.2	0.9
sat	0.3	0.4	0.6	0.1
on	0.2	0.8	0.1	0.5
mat	0.4	0.3	0.7	0.2

Step 2: CNN filters slide over embeddings (captures context)

Step 3: For each token, output probabilities:

Token	PERSON	ORG	GPE	DATE	O
The	0.01	0.01	0.01	0.02	0.95
cat	0.10	0.02	0.01	0.02	0.85
...					

5 Using SpaCy for NER

Overview

SpaCy provides pre-trained NER models that are easy to use and highly accurate for common entity types.

5.1 Basic SpaCy NER

```

1 import spacy
2 from spacy import displacy
3
4 # Load pre-trained model
5 nlp = spacy.load("en_core_web_sm")
6
7 # Process text
8 text = """Donald Trump won more than 50 electoral votes this year.
9 Tesla's stock rose 2.4% after the Federal Reserve announcement."""
10
11 doc = nlp(text)
12

```

```

13 # Extract entities
14 for ent in doc.ents:
15     print(f"{ent.text:20} {ent.label_[:10]} {ent.start_char}-{ent.end_char}")

```

Listing 5: SpaCy NER Basics

Output:

Donald Trump	PERSON	0-12
more than 50	CARDINAL	17-29
this year	DATE	47-56
Tesla	ORG	58-63
2.4%	PERCENT	78-82
Federal Reserve	ORG	94-109

5.2 Visualizing Entities

```

1 # Render in Jupyter notebook
2 displacy.render(doc, style="ent", jupyter=True)
3
4 # Or save to HTML file
5 html = displacy.render(doc, style="ent", page=True)
6 with open("ner_visualization.html", "w") as f:
7     f.write(html)

```

Listing 6: Visualize NER with displacy

5.3 SpaCy NER Label Reference

Label	Description
PERSON	People, including fictional
NORP	Nationalities, religious, political groups
FAC	Facilities (buildings, airports, highways)
ORG	Companies, agencies, institutions
GPE	Countries, cities, states
LOC	Non-GPE locations
PRODUCT	Objects, vehicles, foods
EVENT	Named hurricanes, battles, wars
WORK_OF_ART	Titles of books, songs
LAW	Named documents made into laws
DATE	Absolute or relative dates
TIME	Times smaller than a day
PERCENT	Percentage
MONEY	Monetary values
QUANTITY	Measurements
ORDINAL	“first”, “second”, etc.
CARDINAL	Numerals not falling into other categories

6 Fine-Tuning SpaCy NER

Overview

Pre-trained NER models may not recognize domain-specific entities. SpaCy allows fine-tuning on custom data to add new entity types or improve accuracy.

6.1 When to Fine-Tune

Important

Consider fine-tuning when:

- Domain-specific entities (drug names, product codes)
- New entity categories not in default model
- Improving accuracy for your specific text type
- Handling industry jargon or technical terms

6.2 Training Data Format

```

1 # Training data format: (text, {"entities": [(start, end, label)]})
2 train_data = [
3     ("Cars in China are selling well",
4      {"entities": [(0, 4, "VEHICLE")]}) ,
5
6     ("Tesla has a lot on the line as an electric vehicle maker",
7      {"entities": [(0, 5, "ORG"), (40, 56, "VEHICLE")]}) ,
8
9     ("My family loves our Honda Civic",
10      {"entities": [(23, 34, "VEHICLE")]}) ,
11
12     ("This car is the best",
13      {"entities": [(5, 8, "VEHICLE")]}) )
14 ]

```

Listing 7: SpaCy Training Data Format

6.3 Fine-Tuning Process

```

1 import spacy
2 from spacy.training import Example
3 import random
4
5 # Load existing model
6 nlp = spacy.load("en_core_web_sm")
7
8 # Get the NER component
9 ner = nlp.get_pipe("ner")
10
11 # Add new entity label
12 ner.add_label("VEHICLE")
13
14 # Disable other pipes during training
15 other_pipes = [pipe for pipe in nlp.pipe_names if pipe != "ner"]
16
17 # Training loop

```

```

18 with nlp.disable_pipes(*other_pipes):
19     optimizer = nlp.resume_training()
20
21     for iteration in range(20):
22         random.shuffle(train_data)
23         losses = []
24
25         for text, annotations in train_data:
26             doc = nlp.make_doc(text)
27             example = Example.from_dict(doc, annotations)
28             nlp.update([example], drop=0.5, losses=losses)
29
30         print(f"Iteration {iteration}, Losses: {losses}")
31
32 # Test the model
33 doc = nlp("I bought a Toyota Camry yesterday")
34 for ent in doc.ents:
35     print(f"{ent.text}: {ent.label_}")

```

Listing 8: Fine-Tuning SpaCy NER

Fine-Tuning Pitfalls

- **Catastrophic forgetting:** Model may “forget” original entities
- **Insufficient data:** Need many examples per entity type
- **Label consistency:** Annotations must be consistent
- Always include some original data in training to prevent forgetting

7 Part-of-Speech Tagging

Overview

Part-of-Speech (POS) tagging identifies grammatical categories (noun, verb, adjective, etc.) for each word. It’s closely related to NER and often used as a preprocessing step.

7.1 Common POS Tags

Tag	Description	Example
NN	Noun, singular	cat, dog, house
NNS	Noun, plural	cats, dogs, houses
NNP	Proper noun, singular	John, London
VB	Verb, base form	run, eat, be
VBD	Verb, past tense	ran, ate, was
VBG	Verb, gerund	running, eating
JJ	Adjective	big, red, beautiful
RB	Adverb	quickly, very, well
IN	Preposition	in, on, at, by
DT	Determiner	the, a, an
CC	Coordinating conjunction	and, but, or
TO	“to”	to (as in “to run”)

7.2 POS Tagging with NLTK

```

1 import nltk
2 from nltk import word_tokenize, pos_tag
3
4 # Download required data
5 nltk.download('punkt')
6 nltk.download('averaged_perceptron_tagger')
7 nltk.download('tagsets')
8
9 # View tag descriptions
10 nltk.help.upenn_tagset('NN')    # Noun
11 nltk.help.upenn_tagset('VB')    # Verb
12
13 # POS tagging
14 text = "The quick brown fox jumps over the lazy dog"
15 tokens = word_tokenize(text)
16 pos_tags = pos_tag(tokens)
17
18 print(pos_tags)
19 # [('The', 'DT'), ('quick', 'JJ'), ('brown', 'JJ'),
20 # ('fox', 'NN'), ('jumps', 'VBZ'), ('over', 'IN'),
21 # ('the', 'DT'), ('lazy', 'JJ'), ('dog', 'NN')]

```

Listing 9: POS Tagging with NLTK

7.3 POS Tagging with SpaCy

```

1 import spacy
2
3 nlp = spacy.load("en_core_web_sm")
4 doc = nlp("The quick brown fox jumps over the lazy dog")
5
6 for token in doc:
7     print(f"{token.text}:10 {token.pos_:6} {token.tag_}")

```

Listing 10: POS Tagging with SpaCy

8 Enhancing Classification with NER

Overview

NER features can improve text classification by providing structured information about document content.

8.1 Feature Engineering with NER

NER-Based Features

Count-based features:

- Number of persons mentioned
- Number of organizations
- Number of locations
- Number of dates/times

Binary indicators:

- Contains person name? (0/1)
- Contains organization? (0/1)
- Contains specific entity (e.g., “Tesla”)? (0/1)

8.2 Implementation Example

```

1 import spacy
2 import pandas as pd
3 from sklearn.feature_extraction.text import TfidfVectorizer
4 from sklearn.model_selection import train_test_split
5 from sklearn.neural_network import MLPClassifier
6 import numpy as np
7
8 nlp = spacy.load("en_core_web_sm")
9
10 def extract_ner_features(text):
11     """Extract NER-based features from text"""
12     doc = nlp(text)
13
14     features = {
15         'n_persons': 0,
16         'n_orgs': 0,
17         'n_gpes': 0,
18         'n_dates': 0,
19         'n_money': 0,
20         'n_percent': 0
21     }
22
23     for ent in doc.ents:
24         if ent.label_ == 'PERSON':
25             features['n_persons'] += 1
26         elif ent.label_ == 'ORG':
27             features['n_orgs'] += 1
28         elif ent.label_ == 'GPE':
29             features['n_gpes'] += 1
30         elif ent.label_ == 'DATE':
31             features['n_dates'] += 1
32         elif ent.label_ == 'MONEY':
33             features['n_money'] += 1
34         elif ent.label_ == 'PERCENT':
35             features['n_percent'] += 1
36
37     return features
38
39 # Extract NER features for all documents
40 ner_features = [extract_ner_features(text) for text in documents]
41 ner_df = pd.DataFrame(ner_features)
42
43 # Combine TF-IDF and NER features
44 tfidf = TfidfVectorizer(max_features=20)
45 X_tfidf = tfidf.fit_transform(documents).toarray()
46 X_combined = np.hstack([X_tfidf, ner_df.values])
47
48 # Train classifier
49 X_train, X_test, y_train, y_test = train_test_split(
50     X_combined, labels, test_size=0.2, random_state=42
51 )
52
53 clf = MLPClassifier(hidden_layer_sizes=(50,), max_iter=500)

```

```

54 clf.fit(X_train, y_train)
55 accuracy = clf.score(X_test, y_test)
56 print(f"Accuracy with NER features: {accuracy:.4f}")

```

Listing 11: NER Feature Enhancement

8.3 Entity-Based Binary Features

```

1 def get_entity_dummies(documents):
2     """Create binary indicators for each unique entity"""
3     all_entities = set()
4
5     # First pass: collect all unique entities
6     for text in documents:
7         doc = nlp(text)
8         for ent in doc.ents:
9             all_entities.add((ent.text, ent.label_))
10
11    # Second pass: create binary features
12    feature_matrix = []
13    for text in documents:
14        doc = nlp(text)
15        doc_entities = set((ent.text, ent.label_) for ent in doc.ents)
16
17        row = [1 if entity in doc_entities else 0
18               for entity in all_entities]
19        feature_matrix.append(row)
20
21    columns = [f"{text}_{label}" for text, label in all_entities]
22    return pd.DataFrame(feature_matrix, columns=columns)
23
24 entity_features = get_entity_dummies(documents)

```

Listing 12: Binary Entity Indicators

9 One-Page Summary

Summary

Named Entity Recognition (NER) identifies and classifies named entities in text.
Common Entity Types:

- PERSON, ORG, GPE, LOC, DATE, TIME, MONEY, PERCENT, CARDINAL

Two Approaches:

Rule-Based:

- Uses regular expressions and dictionaries
- No training data needed
- Limited by predefined patterns
- No context awareness

Statistical/Neural:

- Learns from labeled data
- Uses context for disambiguation

- CNNs capture local context via filters
- SpaCy: `nlp = spacy.load("en_core_web_sm")`

NER as Sequence Labeling:

- BIO scheme: B-TYPE (begin), I-TYPE (inside), O (outside)
- Each token gets a tag
- Output: softmax probabilities over entity types

SpaCy Usage:

```
1 import spacy
2 nlp = spacy.load("en_core_web_sm")
3 doc = nlp("Donald Trump visited Tesla.")
4 for ent in doc.ents:
5     print(ent.text, ent.label_)
```

Fine-Tuning: Add custom entity types with labeled examples. Watch for catastrophic forgetting.

Feature Enhancement:

- Add entity counts to feature vectors
- Create binary entity indicators
- Combine with TF-IDF for classification

Applications: Translation, information extraction, question answering, sentiment analysis, search.

10 Glossary

Key Terms

- **NER:** Named Entity Recognition—identifying entities in text
- **Named Entity:** Real-world object with a name (person, organization, place)
- **BIO Tagging:** Begin-Inside-Outside scheme for sequence labeling
- **POS Tagging:** Part-of-Speech tagging—grammatical categories
- **Regular Expression:** Pattern matching syntax for text
- **Rule-Based NER:** Pattern-matching approach to entity extraction
- **Statistical NER:** Machine learning approach to entity extraction
- **HMM:** Hidden Markov Model—probabilistic sequence model
- **CRF:** Conditional Random Fields—discriminative sequence model
- **SpaCy:** Industrial-strength NLP library for Python
- **Fine-Tuning:** Continuing training on domain-specific data

- **Catastrophic Forgetting:** Model losing original knowledge during fine-tuning
- **GPE:** Geopolitical Entity (countries, cities, states)
- **NORP:** Nationalities, religious, or political groups
- **displacy:** SpaCy's visualization module
- **NLTK:** Natural Language Toolkit—Python NLP library
- **Context:** Surrounding words that help disambiguate meaning
- **Word Boundary:** \b in regex—edges of words