

Word Sense Disambiguation(WSD)

Word Sense Disambiguation (WSD) is the task of determining the correct meaning (sense) of a word based on its context, when the word has multiple possible meanings.

Why is WSD Important in NLP?

Many words in natural language are **ambiguous** — they have more than one meaning. Choosing the correct sense is essential for:

- Accurate understanding
- Machine translation
- Information retrieval
- Speech recognition
- Chatbot interactions

Example of Word Sense Ambiguity:

Consider the word “**bank**”:

- *She sat on the **bank** of the river.* → **bank** = riverbank (land)
- *He deposited money in the **bank**.* → **bank** = financial institution

WSD helps identify the correct sense based on **contextual clues**.

Purpose of WSD:

- To improve semantic understanding in NLP tasks.
- To enable correct translations in machine translation.
- To support question answering systems in finding accurate answers.
- To enhance information extraction from text.

Types of WSD Approaches:

Approach	Description	Example Method
Knowledge-Based	Uses dictionaries, thesauri, WordNet, and semantic rules.	Lesk Algorithm
Supervised Learning	Requires labeled training data (sense-annotated corpus).	Decision Trees, SVM
Unsupervised Learning	Clusters word usage from untagged data using context similarity.	Clustering, Word Embeddings
Semi-Supervised / Bootstrapping	Uses a small set of labeled data to learn patterns and expand.	Yarowsky's Algorithm
Deep Learning / Neural Models	Uses contextual embeddings to learn senses.	BERT, ELMo, Transformer-based models

Real-Time Examples of WSD Applications:

a. Machine Translation:

- Sentence: *"I went to the bank to fish."*
- Without WSD: Translates **bank** as financial bank.
- With WSD: Translates **bank** as riverbank in the target language.

b. Voice Assistants / Chatbots:

- User: *"What's the interest rate at the bank?"*
- Bot needs to know **bank = financial institution**, not riverbank, to respond correctly.

c. Search Engines:

- Query: *"Apple benefits"*
- WSD decides if **Apple = fruit** or **Apple = tech company** based on user history/context.

d. Medical NLP:

- Term: “Cold”
 - Can mean low temperature or illness.
 - WSD is vital in extracting correct patient symptoms or conditions.

LESK ALGORITHM

Example Using Lesk Algorithm (Knowledge-Based WSD):dictionaries-thesauri

Sentence: “He went to the bank to deposit money.”

Senses of "bank" (from WordNet):

- Financial institution
- Land beside a river
- Sloping mound of earth

Context words: "deposit", "money"

Lesk Method:

Compares dictionary definitions of each sense with the words in the sentence.

- Sense 1 (financial institution) has overlap with “money”, “deposit”
- Sense 2 (riverbank) does not
→ **Sense 1 is selected.**

Importance of WSD in NLP:

Benefit	Explanation
Improves Translation	Ensures correct sense is translated.
Enhances Chatbot Responses	Prevents misunderstandings.
Boosts Search Relevance	Filters ambiguous queries properly.
Aids Semantic Parsing	Enables deeper sentence understanding.

Crucial for WSD-Based Applications

Sentiment analysis, summarization, entity linking, etc.

Challenges in WSD:

- Scarcity of large, sense-annotated corpora.
- Contexts are often vague or short.
- New word senses emerge over time.
- Domain adaptation (e.g., "virus" in medicine vs. computer science).

Word Sense Disambiguation (WSD) – Methods (Supervised, Dictionary-Based, Bootstrapping)

1. Supervised Learning Approach

Definition:

Supervised WSD uses sense-annotated corpora (text where word senses are already labeled) to train a machine learning model that can predict the correct sense of a word based on its surrounding context.

How it works:

1. Input: Sense-annotated dataset (e.g., SemCor)
2. Features Extracted: Surrounding words, part-of-speech, syntactic structure, etc.
3. Model Trained: Classifier (e.g., Decision Trees, SVMs, Neural Networks)
4. Prediction: Model assigns the correct sense for ambiguous words in unseen texts.

Example:

- Word: "Java"

- Context 1: *"I love drinking hot Java in the morning."* → Sense = Coffee
- Context 2: *"I am debugging Java code."* → Sense = Programming Language
- A supervised model learns to disambiguate based on context and labeled examples.

Pros:

- High accuracy (when good data is available)
- Effective in domain-specific tasks

Cons:

- Requires large sense-tagged corpora
- Limited scalability to new languages/domains

2. Dictionary and Thesaurus-Based Approach (Knowledge-Based)

Definition:

This method uses existing lexical resources such as WordNet, dictionaries, or thesauri to infer the correct sense of a word by comparing word definitions (glosses) and semantic relations.

Lesk Algorithm (1986)

- Compares glosses of the target word's senses with the words in its surrounding context.
- The sense with the most overlapping words is selected.

Example:

- **Sentence:** *"He deposited cash at the bank."*
- **Senses of "bank":**
 1. Financial institution → Gloss includes "money, deposit"

2. Riverbank → Gloss includes “land, river”

- Overlap with “cash” and “deposit” → **Select financial institution**

Pros:

- No need for labeled data
- Easy to apply across languages (with multilingual WordNets)

Cons:

- Limited accuracy
- Doesn't handle context variability well
- Gloss definitions may not be enough

3. Bootstrapping (Semi-Supervised Learning)

Definition:

Bootstrapping starts with a small set of manually labeled examples (seeds) and learns to identify patterns to automatically label more data iteratively.

Yarowsky Algorithm (1995)

- Assumes:
 - One sense per collocation (words around a target word usually indicate the same sense)
 - One sense per discourse (a word retains the same meaning throughout a document)

How it works:

1. Manually tag a few examples of each sense
2. Train classifier using features like surrounding words

3. Apply classifier to unlabeled data
4. Add high-confidence predictions to training data
5. Repeat until convergence

Example:

- **Word:** "Plant"
 - Seed 1: "The **plant** grows in the garden." → sense = living organism
 - Seed 2: "The **plant** produces steel." → sense = industrial factory
 - Contextual patterns like "grows" vs. "produces" help expand the sense classifier.

Pros:

- Requires less labeled data than supervised
- Adaptable to new domains

Cons:

- Errors can propagate
- Initial seeds must be chosen carefully

Comparison Table:

Method	Data Needed	Accuracy	Scalability	Key Tools/Examples
Supervised	Labeled corpus	High	Low	SemCor, Senseval, SVM, BERT
Dictionary-Based	WordNet, glosses	Medium	High	Lesk Algorithm, WordNet
Bootstrapping	Small seed examples	Medium-High	Medium	Yarowsky Algorithm, Clustering

Python code examples for the following **Word Sense Disambiguation (WSD)** methods:

1. **Lesk Algorithm (Dictionary-Based)** using **NLTK** and **WordNet**
2. **Supervised WSD (Simple Example using Scikit-learn)**

1. Lesk Algorithm in Python using NLTK

```
python
CopyEdit
# Install required libraries first
# pip install nltk

import nltk
from nltk.corpus import wordnet as wn
from nltk.wsd import lesk
from nltk.tokenize import word_tokenize

# Download necessary NLTK data
nltk.download('wordnet')
nltk.download('omw-1.4')
nltk.download('punkt')

# Example sentence
sentence = "He went to the bank to deposit money."
words = word_tokenize(sentence)

# Apply Lesk Algorithm
sense = lesk(words, 'bank') # 'bank' is the target ambiguous word

print("Predicted Sense:", sense)
print("Sense Definition:", sense.definition())
```

Output:


```
Predicted Sense: Synset('depository_financial_institution.n.01')
Sense Definition: a financial institution that accepts deposits
```

This shows Lesk correctly picked **"financial bank"** based on words like “deposit” and “money.”

2. Supervised WSD (Simplified Example using Scikit-learn)

Let's simulate a tiny dataset for the word **"bass"** with senses: `fish` and `instrument`.

```
from sklearn.feature_extraction.text import CountVectorizer
from sklearn.naive_bayes import MultinomialNB
```

```
# Training data
```

```
X_train = [
    "I caught a bass in the lake",
    "The bass was swimming in the pond",
    "He played the bass in the orchestra",
    "She tuned the bass guitar"
]
y_train = ["fish", "fish", "instrument", "instrument"]
```

```
# Convert text to feature vectors
```

```
vectorizer = CountVectorizer()
X_train_counts = vectorizer.fit_transform(X_train)
```

```
# Train a Naive Bayes classifier
```

```
clf = MultinomialNB()
clf.fit(X_train_counts, y_train)
```

```
# Test sentence
```

```
test_sentence = ["He caught a bass with a fishing rod"]
X_test = vectorizer.transform(test_sentence)
```

```
# Predict
```

```
prediction = clf.predict(X_test)
print("Predicted Sense:", prediction[0])
```

Output:

```
Predicted Sense: fish
```

This shows the model correctly classifies “**bass**” as a **fish** based on context words like "caught" and "fishing rod."

Bootstrapping WSD Example (Semi-Supervised Learning)

We'll simulate the **Yarowsky bootstrapping** idea using a small labeled seed set and unlabeled data. classify "**plant**" (as **factory** or **living thing**).

(simple pattern bootstrapping using keyword matching.)

Step-by-step Code:

```
# Seed examples (labeled)
seed_data = [
    ("The plant grows in the garden", "living"),
    ("The chemical plant exploded last night", "factory")
]

# Extract keywords for each sense
living_keywords = {"grows", "garden", "leaf", "watered"}
factory_keywords = {"chemical", "exploded", "build", "machines",
                    "workers"}

# Unlabeled data to classify
unlabeled_data = [
    "The plant was watered daily",
```

```

    "Many workers were at the plant",
    "The plant had green leaves",
    "They built a new plant downtown",
    "He visited the plant during the tour"
]

# Classification function based on keyword overlap
def classify(sentence):
    tokens = set(sentence.lower().split())
    living_score = len(tokens & living_keywords)
    factory_score = len(tokens & factory_keywords)

    if living_score > factory_score:
        return "living"
    elif factory_score > living_score:
        return "factory"
    else:
        return "unknown"

# Classify each sentence
for sentence in unlabeled_data:
    label = classify(sentence)
    print(f"Sentence: {sentence}\n→ Predicted Sense: {label}\n")

```

Output:

```

Sentence: The plant was watered daily
→ Predicted Sense: living

```

```

Sentence: Many workers were at the plant
→ Predicted Sense: factory

```

```

Sentence: The plant had green leaves
→ Predicted Sense: living

```

```

Sentence: They built a new plant downtown

```

→ Predicted Sense: factory

Sentence: He visited the plant during the tour

→ Predicted Sense: unknown