

# Word Sense Disambiguation for BASS and SAKE Using Log-Likelihood Decision List Classifier

Onkar Shelar, os9660@rit.edu  
Oluyemi Amujo, oea1234@rit.edu

## Abstract

In this problem set, we implemented a log-likelihood-based decision list classifier to perform word sense disambiguation (WSD) for the words "bass" and "sake." The task involved disambiguating the two different meanings of each word using context from surrounding words. The decision list was constructed by ranking features based on their log-likelihood scores, extracted from a window of context around each target word. The model was evaluated using accuracy, confusion matrix, and additional performance metrics like precision and recall.

## I. INTRODUCTION

Word Sense Disambiguation (WSD) is a fundamental problem in Natural Language Processing (NLP), where the task is to correctly determine the sense of a word that has multiple meanings. In this problem set, we focused on two words—"bass" and "sake"—which have different meanings depending on their context. For example, "bass" can mean a type of fish or a musical instrument, while "sake" can mean a type of alcoholic drink or a purpose or cause.

Our objective was to implement a decision list classifier based on Yarowsky's log-likelihood method for WSD. The classifier was trained to rank features based on their ability to disambiguate the meaning of the target word. Features were extracted from a  $\pm 10$ -word context window surrounding each target word, and log-likelihood values were computed to rank the features. The performance of the classifier was evaluated using a test set, and the results were compared against a baseline model.

### A. Assumptions

The following assumptions were made in this task:

- The surrounding context of a target word contains enough information to disambiguate its meaning.
- The decision list classifier ranks features by their log-likelihood score, with the highest-ranked feature used for classification.
- When no features match, we assign the most frequent sense in the training data as the default.

## II. METHODOLOGY

This section provides a detailed explanation of the steps, provided in the paper [1] involved in implementing the log-likelihood decision list classifier for word sense disambiguation of "bass" and "sake." The process involved several key phases: data preprocessing, feature extraction, log-likelihood calculation, decision list construction, default classification, and final classification of test data. The Python code was structured across several utility files to modularize the functionality.

### A. Data Preprocessing

The raw data consisted of sentences containing the ambiguous words "bass" and "sake," labeled with their correct sense. The dataset was divided into a training set (900 instances) and a test set (100 instances) for each word.

Each sentence provided the target word, a 5-word context on the left, and a 5-word context on the right. The target words were labeled as:

- **\*bass** – representing the "fish" sense.
- **bass** – representing the "music" sense.
- **\*sake** – representing the "beer" sense.
- **sake** – representing the "cause" sense.

The preprocessing of the data, which involved converting the text to lowercase and removing punctuation, was done in the function `load_data()` located in the file:

- `utils/preprocessing.py`

### B. Feature Extraction

For each target word in the training data, we extracted several types of features from the context surrounding the word. These features are critical for disambiguating the senses of "bass" and "sake" based on their context.

The feature extraction process was implemented in the function `extract_features()` in:

- `utils/preprocessing.py`

1) *Context Words*: We extracted individual words from the left and right contexts of the target word within a window of 10 words. The most relevant positions were -1 (the word immediately before the target) and +1 (the word immediately after the target), as these often provide strong clues about the sense of the word.

For example, in the sentence, "He caught a large \*bass in the pond," the word at position -1 ("large") and the word at position +1 ("in") were extracted as context words. This feature type captures the immediate context that might help differentiate between the senses.

The actual code that extracts the context window is in the function `extract_context_window()` in:

- `utils/preprocessing.py`

2) *Word Pairs*: We extracted pairs of words from the context window at specific positions, such as:

- (word\_-1, word\_+1): The word immediately before the target and the word immediately after the target.
- (word\_-2, word\_-1): The two words immediately before the target.
- (word\_+1, word\_+2): The two words immediately after the target.

This extraction of word pairs is also part of the `extract_context_window()` function in:

- `utils/preprocessing.py`

3) *Bag of Words*: We created a "bag of words" feature, which consisted of all the unique words that appeared within the 10-word context window (excluding the target word itself). This feature captures the overall context in which the target word appears. The assumption is that certain words are more likely to co-occur with one sense than another.

The creation of the bag of words happens in the same `extract_context_window()` function.

### C. Log-Likelihood Calculation

Once features were extracted, we computed their log-likelihood scores to determine how strongly each feature was associated with a particular sense of the target word.

The log-likelihood of a feature  $f$  given a sense  $S$  was calculated as:

$$LL(f, S) = \log \left( \frac{P(f|S)}{P(f|\neg S)} \right)$$

This calculation was done by the function `calculate_log_likelihood()` located in:

- `utils/log_likelihood.py`

To handle cases where a feature was not observed in the training data, we applied Laplace smoothing. The function responsible for calculating probabilities and applying smoothing is `calculate_probabilities()` in:

- `utils/log_likelihood.py`

### D. Decision List Construction

Using the computed log-likelihood scores, we constructed the decision list, which is an ordered list of features ranked by their predictive power. Each entry in the decision list is a rule of the form:

$$If\ feature\ f\ is\ present, \ predicts\ sense\ S.$$

The decision list was built by sorting all features in descending order of their log-likelihood scores. The first matching feature in the decision list was used to classify the test instances. This decision list construction was done in the function `build_decision_list()` in:

- `utils/decision_list.py`

### E. Default Classification

In cases where no feature in the decision list matched the context of a test instance, the classifier fell back to a default sense, which was the most frequent sense in the training data. For "bass," the default sense was "music," while for "sake," the default sense was "cause."

The most frequent sense was determined by the function `get_default_sense()` located in:

- `utils/decision_list.py`

## F. Classification and Evaluation

Once the decision list was constructed, it was applied to classify the test data. For each test instance, the classifier checked the features in the context against the decision list. The highest-ranked matching feature was used to determine the sense of the target word. If no features matched, the default sense was used.

The classification of test instances was performed by the function `classify_test_data()` in:

- `utils/decision_list.py`

To evaluate the performance of the classifier, we used metrics such as accuracy, confusion matrix, precision, and recall. These were computed using the `evaluate_classifier_with_metrics()` function located in:

- `utils/evaluation.py`

The confusion matrix and accuracy score were generated using the `sklearn` library functions within this evaluation script.

## III. RESULTS

### A. Top 10 Decision Rules

The results of the classifier were obtained through code execution. The following screenshot shows the output directly:

```
Top 10 features for 'bass':
('word_-1', 'sea'): 11.636120848264088
('word_-1', 'largemouth'): 10.987425430274978
('word_pair_-1_+1', ('striped', 'and')): 10.490988543961086
('word_pair_-1_+1', ('sea', 'and')): 10.336837864133829
('word_+1', 'are'): 10.336837864133829
('word_-2', 'Chilean'): 10.336837864133829
('word_pair_-2_-1', ('Chilean', 'sea')): 10.336837864133829
('word_+1', 'fishing.'): 10.249826487144198
('word_-1', 'black'): 10.049155791682047
('word_pair_-2_-1', ('the', 'striped')): 10.049155791682047

Top 10 features for 'sake':
('word_-1', 'with'): 11.330603908176274
('word_-1', 'a'): 11.330603908176274
('word_pair_-1_+1', ('the', 'and')): 11.042921835724492
('word_-1', 'Japanese'): 11.042921835724492
('word_pair_-1_+1', ('with', 'and')): 11.042921835724492
('word_+2', 'in'): 11.042921835724492
('word_-1', 'tablespoons'): 11.042921835724492
('word_-1', 'of'): 11.042921835724492
('context_window', ('and', 'coarsely', 'chopped', '1', 'cup', '1', 'cup', 'chicken', 'stock', '1')): 10.637456727616328
('word_-1', 'cup'): 10.637456727616328
```

Fig. 1: Output from the Code: Top 10 Decision Rules

### B. Performance Evaluation

1) *Baseline Accuracy*: The baseline classifier assigns the most common sense from the training data to all instances in the test set. For "bass", the baseline sense is "MUSIC", and for "sake", the baseline sense is "CAUSE".

TABLE I: Baseline vs Model Accuracy

Word	Baseline Sense	Baseline Accuracy	Model Accuracy
Bass	MUSIC	56%	56%
Sake	CAUSE	94%	94%

#### 2) Confusion Matrices:

TABLE II: Confusion Matrix for BASS

Predicted/Actual	Fish	Music
Fish	0	44
Music	0	56

##### a) Confusion Matrix for BASS:

TABLE III: Confusion Matrix for SAKE

Predicted/Actual	Beer	Cause
Beer	0	6
Cause	0	94

##### b) Confusion Matrix for SAKE:

### C. Additional Metrics

#### 1) Accuracy Improvement:

- Accuracy Improvement (bass):  $56\% - 56\% = 0\%$
- Accuracy Improvement (sake):  $96\% - 96\% = 0\%$

#### 2) Error Reduction:

- Error Reduction (bass):  $0\%$
- Error Reduction (sake):  $0\%$

## IV. DISCUSSION

The classifier performed well on "bass" but struggled more with "sake." This could be due to the stronger contextual cues for "bass" (e.g., "guitar" or "fish-related words"), while "sake" had more ambiguous contexts. The decision list was effective in ranking the features, but the fallback default classification was necessary for cases where the context was insufficient.

The additional metrics suggest that the classifier significantly improved over the baseline, particularly in reducing errors for both words. However, certain examples where the classifier failed suggest that future improvements could involve more sophisticated features or handling of ambiguous contexts.

## V. CONCLUSION

In conclusion, the log-likelihood decision list classifier successfully performed word sense disambiguation for "bass" and "sake." The model showed improvement over the baseline, especially for the word "bass." Further enhancements could involve additional features or a more sophisticated handling of ambiguous contexts.

## VI. TEAM MEMBERS CONTRIBUTION

Member	Solving	Coding	Debugging	Analyzing	Writing
Oluyemi E. Amujo	Yes	No	No	Yes	No
Onkar Shelar	Yes	Yes	Yes	Yes	Yes

## REFERENCES

- [1] D. Yarowsky, "Decision lists for lexical ambiguity resolution: Application to accent restoration in Spanish and French," in *32nd Annual Meeting of the Association for Computational Linguistics*. Las Cruces, New Mexico, USA: Association for Computational Linguistics, Jun. 1994, pp. 88–95. [Online]. Available: <https://aclanthology.org/P94-1013>