# 1. Objective

To disambiguate word pairs using a Naive-Bayesian technique and answer to some questions.

# 2. Installation

**Programming language:** `Python`

**Source code location:** `/home/sbillah/nlp2/`

**Corpus Selection:** `I select "`**`senseval"`**` corpus, which is specially designed for WSD. Here is the download link:` **`http://www.senseval.org/.`**

# 3. Design (Naive Bayesian Approach)

## 3.1    Algorithm:

I use two in-memory Hash-tables to store conditional probabilities of contexts associated with the pseudowords.

```
// preprocess corpus files.
for each sense-file f:
     1. replace the word with pseudoword and remember its sense in
        <tag> region.

     2. extract the context words in both side of the pseudoword, and
     store the contexts in two files: training, and testing by 8:2
     ratio.

//training
For each training file f:
     for each line in f:
          1. update C(context-words), C(word), & sense sk in respected
             hash-tables.
     from the counts, compute P(c_i|s_k), p(s_k)and store in hash-tables.

//testing
for each testing file f:
     for each line in f:
          1. apply Laplace smoothing on conditional probabilities.
          2. compute argmax score(s_k) using the formula in the book.
          3. compare the predicted value with actual value.

return accuracy in percentage.
```

**Time Complexity:**

- **Preprocessing phase:** O(# of lines containing word1) + O(#of lines containing word2)
- **Training phase:** 2*O(2*context_size * #lines in training file)
- **Testing phase:** O(2*context_size * #lines in testing file)
- **Overall:** O(2*context_size*(#of word1+ #of word2)

**Space Complexity:**

- **Overall:** O(2 * context_size * (#of word1 +  #of word2) )

## 3.2    Corpus Description:

The **Senseval** WSD corpus has total 35 sense-tagged words. Each word has more than 5 senses. But due to the simplified requirement of our homework, I ignore all those senses. Therefore, for a word pair, I consider only two senses (0,1). Here are my selected word-pairs and their individual occurrence in the corpus.

| Pair | Words | Word Counts |
|---|---|---|
| 1 | amaze | 319 |
|   | behaviour | 1003 |
| 2 | sack | 296 |
|   | sanciton | 101 |
| 3 | knee | 477 |
|   | onion | 29 |
| 4 | accident | 1303 |
|   | wooden | 370 |

Below is a snapshot of some lines from "***accident.cor***"  file (context for word, 'accident'):

```
800001
Late on Thursday night it was travelling at about three metres a second in
wind blowing at 20 to 25 knots when an empty car fell off just as it reached
the top.
The <tag "532675">accident</> appeared to have little effect on the Christmas
party, except to lengthen it considerably.

800002
An image of earnest Greenery is almost tangible.
Eighteen years ago she lost one of her six children in an <tag
"532675">accident</> on Stratford Road, a tragedy which has become a pawn in
the pitiless point-scoring of small-town vindictiveness.
```

```
800003
It's a sentiment I recommend to you all.
The <tag "532675">accident</> occurred on the Saturday of the annual Popular
Flying Association (PFA) rally at Cranfield.

...
```

## 3.3    Context Selection:

I varied context length from 1 to 19 (on both side) as shown in the figure below:



Different level of accuracy is obtained under different context size. The results are given in the next chapter.

## 3.4    Laplace Smoothing:

During testing phase, some context-words are not seen before in training phase. Instead of assigning zero probability for them, I use Laplace Smoothing. The Laplace Smoothing is given below:
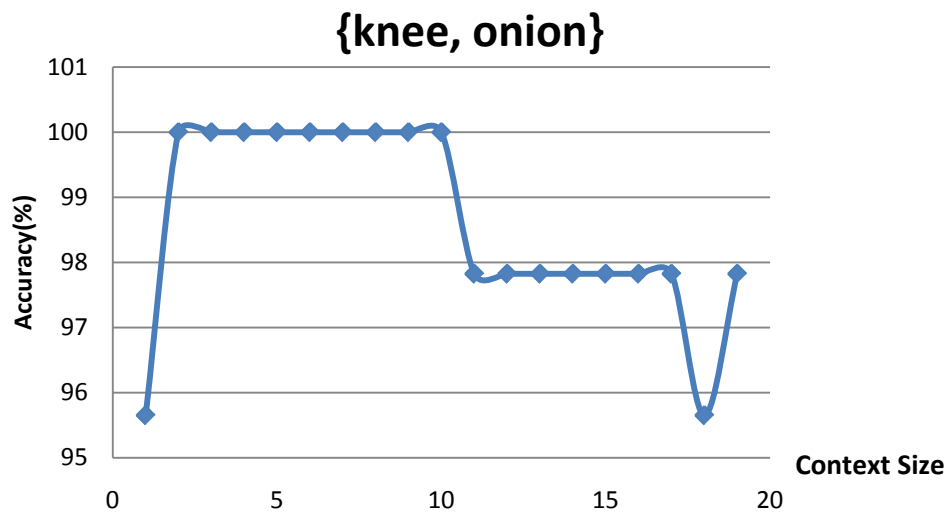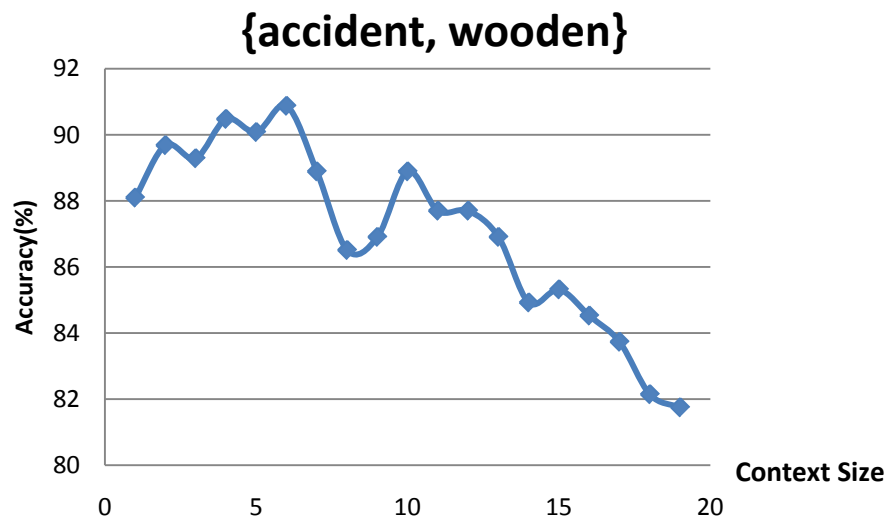
- If $P(ci|s_1) > 0$:
    o  $P(ci|s_1)' = (P(ci|s_1)*10000+1)/( 10000 + context\_size )$
- Else:
    o  $P(ci|s_1)' = 1.0/( 10000 + context\_size )$

# 4  Experimental Result

Here, I provide the experimental results for different context size and word-pairs. After the table, I also present these data graphically for better understanding.

| Context Size | Pair 1 accuracy {accident, wooden} | Pair 2 accuracy {knee, onion} | Pair 3 accuracy {sack, sanction} | Pair 4 accuracy {amaze, behaviour} |
|---|---|---|---|---|
| 1 | 88.09524 | 95.65217 | 85.48387 | 88.04781 |
| 2 | 89.68254 | 100 | 85.48387 | 88.04781 |
| 3 | 89.28571 | 100 | 77.41935 | 83.66534 |
| 4 | 90.47619 | 100 | 77.41935 | 82.47012 |
| 5 | 90.07937 | 100 | 79.03226 | 82.07171 |
| 6 | 90.87302 | 100 | 79.03226 | 81.67331 |
| 7 | 88.88889 | 100 | 85.48387 | 83.26693 |
| 8 | 86.50794 | 100 | 87.09677 | 82.86853 |
| 9 | 86.90476 | 100 | 82.25806 | 81.2749 |

| 10 | 88.88889 | 100 | 82.25806 | 80.47809 |
|----|----------|-----|----------|----------|
| 11 | 87.69841 | 97.82609 | 82.25806 | 80.07968 |
| 12 | 87.69841 | 97.82609 | 80.64516 | 80.87649 |
| 13 | 86.90476 | 97.82609 | 82.25806 | 79.68127 |
| 14 | 84.92063 | 97.82609 | 82.25806 | 80.87649 |
| 15 | 85.31746 | 97.82609 | 80.64516 | 78.88446 |
| 16 | 84.52381 | 97.82609 | 80.64516 | 78.88446 |
| 17 | 83.73016 | 97.82609 | 79.03226 | 77.68924 |
| 18 | 82.14286 | 95.65217 | 79.03226 | 77.29084 |
| 19 | 81.74603 | 97.82609 | 75.80645 | 77.68924 |

## {accident, wooden}



## {knee, onion}

## {sack, sanction}



## {amaze, behaviour}

## Discussion:

1. Accuracy always decrease with the increase of context size.
2. For different word-pairs, maximum accuracy is obtained in different context size.
3. The overall performance of Naive Bayesian disambiguation is around 90%.

# Q&A

2.9 Relative frequency:

**File:** DavidBowie.html
**Entropy:** 3.52728218653

| | | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 2370 | 8 | 917 | h | 7057 | s | 14664 |
| 0 | 1909 | a | 19075 | k | 3306 | r | 12658 |
| 3 | 885 | c | 8595 | j | 423 | u | 4556 |
| 2 | 1502 | b | 4576 | m | 4891 | t | 16047 |
| 5 | 784 | e | 23268 | l | 11101 | w | 5067 |
| 4 | 750 | d | 8301 | o | 11614 | v | 2178 |
| 7 | 881 | g | 3754 | n | 13074 | y | 2469 |
| 6 | 721 | f | 5294 | q | 144 | x | 1038 |
| 9 | 1358 | i | 20024 | p | 7083 | z | 232 |

**File:** Genghis Khan - Wikipedia, the free encyclopedia.html
**Entropy:** 3.55870437955

| | | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 1844 | 8 | 594 | h | 9353 | s | 13442 |
| 0 | 1870 | a | 20479 | k | 4675 | r | 11698 |
| 3 | 873 | c | 7362 | j | 848 | u | 4415 |
| 2 | 1562 | b | 3328 | m | 5568 | t | 15763 |
| 5 | 686 | e | 21399 | l | 11457 | w | 4265 |
| 4 | 566 | d | 7111 | o | 11204 | v | 1788 |
| 7 | 505 | g | 6384 | n | 14745 | y | 2359 |
| 6 | 610 | f | 4749 | q | 263 | x | 1218 |
| 9 | 624 | i | 20236 | p | 6122 | z | 392 |

**File:** Steve Jobs - Wikipedia, the free encyclopedia.html
**Entropy:** 3.61251296863

| | | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 7046 | 8 | 1811 | h | 14120 | s | 31064 |
| 0 | 7798 | a | 36239 | k | 4429 | r | 23907 |
| 3 | 1564 | c | 18155 | j | 2588 | u | 6540 |
| 2 | 6815 | b | 7692 | m | 9194 | t | 32956 |
| 5 | 1808 | e | 47051 | l | 22943 | w | 9403 |
| 4 | 1605 | d | 11058 | o | 23755 | v | 5120 |
| 7 | 1667 | g | 6116 | n | 26050 | y | 4506 |
| 6 | 1966 | f | 10656 | q | 179 | x | 3221 |
| 9 | 2120 | i | 31138 | p | 16747 | z | 504 |

| 1 | 6843 | | 8 | 1769 | | H | 21543 | | s | 34129 |
|---|------|---|---|------|---|---|-------|---|---|-------|
| 0 | 5721 | | a | 45163 | | K | 8332 | | r | 35857 |
| 3 | 1761 | | c | 21142 | | J | 1441 | | u | 10682 |
| 2 | 3863 | | b | 9324 | | M | 12329 | | t | 43826 |
| 5 | 2560 | | e | 56370 | | L | 34877 | | w | 12176 |
| 4 | 1703 | | d | 19968 | | O | 29943 | | v | 5852 |
| 7 | 1273 | | g | 9962 | | N | 34098 | | y | 7013 |
| 6 | 1797 | | f | 13635 | | Q | 445 | | x | 3895 |
| 9 | 2795 | | i | 49316 | | P | 16278 | | z | 654 |

**Finally, the total Corpus frequency:**

| 1 | 18103 | | 8 | 5091 | | H | 52073 | | s | 93299 |
|---|-------|---|---|------|---|---|-------|---|---|-------|
| 0 | 17298 | | a | 120956 | | K | 20742 | | r | 84120 |
| 3 | 5083 | | c | 55254 | | J | 5300 | | u | 26193 |
| 2 | 13742 | | b | 24920 | | M | 31982 | | t | 108592 |
| 5 | 5838 | | e | 148088 | | L | 80378 | | w | 30911 |
| 4 | 4624 | | d | 46438 | | O | 76516 | | v | 14938 |
| 7 | 4326 | | g | 26216 | | N | 87967 | | y | 16347 |
| 6 | 5094 | | f | 34334 | | Q | 1031 | | x | 9372 |
| 9 | 6897 | | i | 120714 | | P | 46230 | | z | 1782 |

**2.10 KL Divergence**

| KL-divergence <1,2> | 0.003688 |
|---------------------|----------|
| KL-divergence <2,1> | 0.003688 |
| **KL-divergence <1,3>** | **0.03306** |
| **KL-divergence <3,1>** | **0.03306** |
| KL-divergence <2,3> | 0.028468 |
| KL-divergence <3,2> | 0.028468 |

So, the corpus 1(english1) and corpus 3 (french1) have the highest score. In fact, these two corpus are same and translation of each other, which justifies the result.