

1. Ziel von Wortbedeutungsdesambiguierung ? Welche Input bekommt der Desambigierer und was gibt er als Output/Ergebnis aus?
2. Nenne ein Beispiel von Input and Output.

## 1. Ziel von Wortbedeutungsdesambiguierung ?

> Für ein mehrdeutiges Wort in einem Satz, die richtige Wortbedeutung zu bestimmen.

Welche Input bekommt der Desambigierer und was gibt er als Output/Ergebnis aus?

> Input: ein Text in dem der zu desambiguierendem Wort vorkommt.  
Output: die wahrscheinlichste Wortbedeutung

## 2. Nenne ein Beispiel von Input and Output.

> Input : "Ich setze mich auf die Bank" , das zu desambiguierende Wort = "Bank"  
> Output: Bank = Bank zum Sitzen

## Desambiguierung

Die Anwendung des Naïve Bayes-Klassifikators ist einfach.

Für ein gegebenes mehrdeutiges Wort  $w$  im Kontext der Wörter  $w_1 \dots w_n$

- berechnet man das Produkt  $p(s) \prod_{i=1}^n p(w_i|s)$  für jede Bedeutung  $s$
- und wählt die Bedeutung mit dem höchsten Wert.

3. Wortbedeutungsdesambiguierung ist eine Klassifikation-Aufgabe. Was wollen wir klassifizieren?

4. Was ist die Grundidee von einem statistischen Model für Wortbedeutungsdesambiguierung (in unserer Vorlesung)? ( Welche Informationen über das zu disambigierende Wort werden verwendet um das Ergebnis des Models zu berechnen?)

Antwort : context wörter, Inhaltswörter

keine Antwort: Bag of Words model, alle substantive Wörter ohne Konjunktionen

Welche Annahme macht Naive Bayes Model : Bag of Words model im Kontext. aber die Kontextwörter (Positionen) sind nicht abhängig voneinander

3. Wortbedeutungsdesambiguierung ist eine Klassifikation-Aufgabe. Was wollen wir klassifizieren?

- Classify the meaning of a word from many meanings (multi class classification).
- Given a word with many possible senses, the classifier the word to the correct sense.

4. Was ist die Grundidee von einem statistischen Model für Wortbedeutungsdesambiguierung (in unserer Vorlesung)? (= welche Annahme macht man über die Bedeutung eines Wortes? / Welche Informationen über das zu disambigierende Wort werden verwendet um das Ergebnis des Models zu berechnen?)

- Context words (neighbouring words of ambiguous word) can tell about the word meaning. E.g. if the ambiguous word is mostly tagged with meaning\_1 when it occurs with a particular group of context words, then the model should classify this word to meaning\_1, when it occurs with these context words.

## Formel von einem statistischen Wortbedeutungsdesambiguierer?

$$s^{\wedge} = \operatorname{argmax}_s p(s \mid \text{context words})$$

$s$  = best sense

$s$  = all senses of the word that we want to disambiguate

context words = for example,  $n$  previous and next  $n$  words of the word that we want to disambiguate

# Formel von einem statistischen Wortartdesambiguierer

Wir wollen die wahrscheinlichste Bedeutung  $\hat{s}$  des Wortes  $w$  im Kontext  $C$  bestimmen:

$$\begin{aligned}\hat{s} &= \arg \max_s p(s|C) \\ &= \arg \max_s \frac{p(s, C)}{p(C)} \\ &= \arg \max_s p(s, C)\end{aligned}$$

$p(C)$  ist eine Konstante, die keinen Einfluss auf das Ergebnis der argmax-Operation hat.

$$= \arg \max_s p(s, C)$$

Wie wird  $p(s, C)$  definiert und welche Annahme wurde gemacht?



Wie wird  $p(s, C)$  definiert und welche Annahme wurde gemacht?

Der Kontext  $C$  wird durch eine Menge von Kontextwörtern  $= w_1 \dots w_n$  (z.B. die 100 nächsten Inhaltswörter) repräsentiert.

$$\begin{aligned} p(s, C) &= p(s, w_1 \dots w_n) \\ &= p(s) \prod_{i=1}^n p(w_i | s, w_1 \dots w_{i-1}) \end{aligned}$$

Zur Vereinfachung des Modelles wird angenommen, dass die Kontextwörter statistisch unabhängig sind, wenn die Bedeutung  $s$  gegeben ist ("Bag of Words"-Modell):

$$p(s, C) = p(s) \prod_{i=1}^n p(w_i | s)$$

**Formel für Naïve Bayes-Klassifikator ? (Erklären auch was die Variablen bedeuten)**

$$\hat{s} = \arg \max_s p(s) \prod_{i=1}^n p(w_i | s)$$

**Welcher Teil heißt Apriori-Wk?**

>  $p(s)$

**Ein Beispiel**

$p(s) = p(\text{Bank} = \text{Sitzplatz}) , p(\text{Bank} = \text{Geldinstitut})$

$p(s | w_1, w_2, w_3) = p(s) p(w_1 | s) p(w_2 | s) p(w_3 | s)$

Ich **setze** mich auf die **Bank** und **trinke Cola**

$p(\text{sense1} | w_1, w_2, w_3) = p(\text{sense1}) \prod p(w_1, w_2, w_3 | \text{sense1})$  – nicht richtig

**Product  $i=1$  bis  $n$  von  $p(w_i | s) = p(w_1 | s) * p(w_2 | s) * p(w_3 | s)$**

**Text:** Der **Hahn** am Waschbecken im Badezimmer tropft.

Wir wollen “Hahn” disambigieren zwischen “Hahn=Tier” und “Hahn=Wasserhahn”

Was muss das Model konkret berechnen um  $s^*$  zu bekommen?  
(Was wäre  $s$ , und  $w_i$  ? )

$$\hat{s} = \arg \max_s p(s) \prod_{i=1}^n p(w_i | s)$$

**Text:** Der Hahn am Waschbecken im Badezimmer tropft.

Wir wollen “Hahn” disambigieren zwischen “Hahn=Tier” und “Hahn=Wasserhahn”

Was müssen wir berechnen um  $s^*$  zu bekommen?

$$p(s, C) = p(s) \prod_{i=1}^n p(w_i | s)$$

**Text:** Der Hahn am Waschbecken im Badezimmer tropft.

**Klasse:** Hahn1 (=Wasserhahn)

**Wahrscheinlichkeit:**  $p(\text{Hahn1}, \text{Waschbecken}, \text{Badezimmer}, \text{tropft}) =$

$p(\text{Hahn1}) p(\text{Waschbecken} | \text{Hahn1}) p(\text{Badezimmer} | \text{Hahn1}) p(\text{tropft} | \text{Hahn1})$

$$\hat{s} = \arg \max_s p(s) \prod_{i=1}^n p(w_i | s)$$

Wie werden  $p(s)$  und  $p(w_i | s)$  geschätzt? (Schrieb Formel und erklären wofür stehen die Variablen  $s^*$ ,  $s$ ,  $w_i$ ,  $n$ )

$p(s) = \dots$

$p(w_i | s) = \dots$

Wie sieht die Training-Daten aus ? Welche Annotationen brauchen wir?

## Parameterschätzung

Um einen Wortbedeutungsdesambiguierer für das Wort  $w$  zu erstellen, braucht man ein Korpus, in dem jedes Vorkommen des Wortes  $w$  manuell mit seiner korrekten Bedeutung annotiert wurde.

Man zählt die Bedeutungen und berechnet die Apriori-Wahrscheinlichkeit:

$$p(s) = \frac{f(s)}{\sum_{s'} f(s')}$$

Dann zählt man alle Wörter  $w'$ , die im Kontext der Bedeutung  $s$  des Wortes  $w$  auftauchen und berechnet:

$$p(w'|s) = \frac{f(w', s)}{\sum_{w''} f(w'', s)}$$

Training data with annotation.

Ich gehe in die **Bank**(Bank= Geldinstitut) und Geld abzuheben.

Ich setze mich auf eine **Bank** (Bank = Sitzplatz) unter einem Baum.

Die Bank ist nicht mehr offen, deswegen kann ich nicht zur **Bank** (Bank= Geldinstitut) gehen.

Eine Frau setzte sich auf die **Bank** (Bank = Sitzplatz) und liest ein Buch.

Training data with annotation.

Ich gehe in die **Bank**(Geldinstitut) und Geld abzuheben.

Ich setze mich auf eine **Bank** (Sitzplatz) unter einem Baum.

Die **Bank** (Geldinstitut) ist nicht mehr offen, deswegen kann ich nicht zur **Bank** (Geldinstitut) gehen.

Eine Frau setzte sich auf die **Bank** (Sitzplatz) und liest ein Buch.

Die **Bank**(Geldinstitut) ist zu.

$$p(s) = \frac{f(s)}{\sum_{s'} f(s')}$$

We want to disambiguate “Bank”. Apply the formula to both senses of Bank and calculate  $f(s)$  and all  $f(s')$



Training data with annotation.

Ich gehe in die **Bank**(Geldinstitut) und Geld abzuheben.

Ich setze mich auf eine **Bank** (Sitzplatz) unter einem Baum.

Die **Bank** (Geldinstitut) ist nicht mehr offen, deswegen kann ich nicht zur **Bank** (Geldinstitut) gehen.

Eine Frau setzte sich auf die **Bank** (Sitzplatz) und liest ein Buch.

Die **Bank**(Geldinstitut) ist zu.

We want to disambiguate “Bank”. Apply the formula to both senses of Bank and calculate  $f(s)$  and all  $f(s')$

$$p(s) = \frac{f(s)}{\sum_{s'} f(s')}$$

$$p(\text{Geldinstitut}) = f(\text{Geldinstitute}) / (f(\text{Geldinstitute}) + f(\text{Sitzplatz}))$$

$$p(\text{Sitzplatz}) = f(\text{Sitzplatz}) / (f(\text{Geldinstitute}) + f(\text{Sitzplatz}))$$

$$f(\text{Geldinstitut}) = 4, f(\text{Sitzplatz}) = 2$$

Training data with annotation.

Ich gehe in die **Bank**(Geldinstitut) und **Geld abzuheben**.

Ich setze mich auf eine **Bank** (Sitzplatz) unter einem Baum.

Die **Bank** (Geldinstitut) ist nicht mehr offen, deswegen kann ich nicht zur **Bank** (Geldinstitut) gehen.

Eine Frau setzte sich auf die **Bank** (Sitzplatz) und liest ein Buch.

Die **Bank**(Geldinstitut) ist zu.

$$p(w'|s) = \frac{f(w', s)}{\sum_{w''} f(w'', s)}$$

New text: “Ich will ein Buch kaufen und zur **Bank** gehen, aber sie ist vermutlich nicht mehr offen”

(context words are all content words in text)

What are  $w'$  ? How many of them do we have here?

Use the formula on this text for all  $w'$  (with  $s =$  Geldinstitute )

Training data with annotation.

Ich **gehe** in die **Bank**(Geldinstitut) und Geld abzuheben.

Ich setze mich auf eine **Bank** (Sitzplatz) unter einem Baum.

Die **Bank** (Geldinstitut) ist nicht mehr offen, deswegen kann ich nicht zur **Bank** (Geldinstitut) gehen.

Eine Frau setzte sich auf die **Bank** (Sitzplatz) und liest ein Buch.

Die **Bank**(Geldinstitut) ist zu, deswegen kann ich kein Geld abheben.

New text: New text: “Ich will ein **Buch kaufen** und zur **Bank gehen**, aber sie ist vermutlich nicht mehr **offen**”

(context words are all content words in text)

What are  $w'$  ? How many of them do we have here?

Use the formula on this text for all  $w'$  (with  $s$  = Geldinstitute )

$w'$  = context words = gehen, offen, Buch, kaufen

Bank = Geldinstitut :  $s_1$

$p(w', s) = p(\text{gehen}, s_1) = f(\text{gehen}, s_1) / f(\text{gehen}, s_1) + f(\text{offen}, s_1) + f(\text{Buch}, s_1) + f(\text{kaufen}, s_1)$

$p(w', s) = p(\text{offen}, s_1) = f(\text{offen}, s_1) / f(\text{gehen}, s_1) + f(\text{offen}, s_1) + f(\text{Buch}, s_1) + f(\text{kaufen}, s_1)$

$p(w', s) = p(\text{Buch}, s_1) = \dots$

$p(w', s) = p(\text{kaufen}, s_1) = \dots$

Training data with annotation.

Ich **gehe** in die **Bank**(Geldinstitut) und Geld abzuheben.

Ich setze mich auf eine **Bank** (Sitzplatz) unter einem Baum.

Die **Bank** (Geldinstitut) ist nicht mehr **offen**, deswegen kann ich nicht zur **Bank** (Geldinstitut) **gehen**.

Eine Frau setzte sich auf die **Bank** (Sitzplatz) und liest ein Buch.

Die **Bank**(Geldinstitut) ist zu, deswegen kann ich kein Geld abheben.

New text: "Ich will ein **Buch kaufen** und zur **Bank gehen**, aber sie ist vermutlich nicht mehr **offen**"

$$p(w'|s) = \frac{f(w', s)}{\sum_{w''} f(w'', s)}$$

f(gehen, Bank=Geldinstitut) = 2

f(offen, Bank=Geldinstitut) = 1

f(Buch, Bank=Geldinstitut) = 0

f(Kaufen, Bank=Geldinstitut) = 0

Assuming we have trained a model. How can we use it to compute the sense of “Schwester” in “Seine einzige Verwandte, die Schwester, liegt im Krankenhaus” ?

- What parameters are stored in the model?
- And which one do we use?
- What probabilities does the model compute? And how to determine the best meaning using this information?

Assuming we have trained a model. How can we use it to compute the sense of “Schwester” in “Seine einzige Verwandte, die Schwester, liegt im Krankenhaus” ?

**- What parameters are stored in the model?**

> all possible senses of the word

>  $p(\text{sense1})$  ,  $p(\text{sense2})$ , ...

>  $p(\text{context\_word}, \text{sense})$  for all possible context found in the training corpus

**- And which one do we use?**

Seine einzige Verwandte, die **Schwester**, liegt im Krankenhaus

**lemmatisierte Kontextwörter:**

einzig, Verwandte, Krankenhaus, liegen

Wahrscheinlichkeiten:

Wort w	$p(w \text{SchwesterV})$	$p(w \text{SchwesterK})$
einzig	0.001	0.00015
Verwandte	0.02	0.0001
liegen	0.0005	0.002
Krankenhaus	0.0001	0.03

$$p(\text{SchwesterV}) = p(\text{SchwesterK}) = 0.5$$

## Beispiel

Seine einzige Verwandte, die **Schwester**, liegt im Krankenhaus

**lemmatisierte Kontextwörter:**

einzig, Verwandte, Krankenhaus, liegen

Wahrscheinlichkeiten:

Wort w	$p(w \text{SchwesterV})$	$p(w \text{SchwesterK})$
einzig	0.001	0.00015
Verwandte	0.02	0.0001
liegen	0.0005	0.002
Krankenhaus	0.0001	0.03

$$p(\text{SchwesterV}) = p(\text{SchwesterK}) = 0.5$$

$$p(\text{SchwesterV}, \text{einzig}, \text{Verwandte}, \text{Krankenhaus}, \text{liegen}) = 5e-13$$

$$p(\text{SchwesterK}, \text{einzig}, \text{Verwandte}, \text{Krankenhaus}, \text{liegen}) = 4.5e-13$$

⇒ Wir wählen hier die Lesart *SchwesterV*.

How does a simple baseline predict word sense?



## Baseline und Obergrenze

Wenn kein Vergleichssystem zur Verfügung steht, bietet sich die folgende Baseline-Methode für Vergleichszwecke an:

**Baseline:** Wähle immer die Lesart, die in den Trainingsdaten am häufigsten war.

Wenn diese Baseline nicht übertroffen wird, (was oft gar nicht so einfach ist), ist das Verfahren nutzlos.



Wenn man wissen möchte, wie viel “Luft” noch für Verbesserungen vorhanden ist, kann man die menschliche Leistung bei dieser Aufgabe als **Obergrenze** zum Vergleich heranziehen.

**Aufgabe 6)** Erklären Sie, wie eine Pseudowort-Evaluierung funktioniert und wofür man Sie verwenden kann. (3 Punkte)

## **Pseudowort-Evaluierung**

- Wozu benutzen wir Pseudowort-Evaluierung?
- Schritte ?

## Pseudowort-Evaluierung

- Wozu benutzen wir Pseudowort-Evaluierung?
  - Um den Ansatz (Naive Bayes Model) zu evaluieren
  - Vorteile: keine annotiertes Korpus nötig
  - Es kann das wirkliche Model nicht evaluieren
- Schritte ?
  - Model ist trainiert, um das wort “w1-w2” mit zwei Bedeutungen w1, w2 zu disambiguieren.

# Pseudowort-Evaluierung

Die manuelle Annotation von Daten ist teuer.

Zur Evaluierung einer Desambiguierungsmethode kann man jedoch ein großes perfekt annotiertes Trainingskorpus erstellen, indem man zwei Wörter zu einem neuen mehrdeutigen "Pseudo"-Wort zusammenfasst:

## **Pseudowort-Evaluierung:**

- Nimm ein großes Textkorpus.
- Wähle zwei Wörter  $w_1$  und  $w_2$  und ersetze alle Vorkommen von  $w_1$  und  $w_2$  durch das Pseudowort  $w_1-w_2$ .
- Behalte das ursprüngliche Wort als „Bedeutungs“-Annotation.
- Teile das Korpus in zwei Teile.
- Trainiere den Desambiguierer auf dem ersten Teil, zwischen den Bedeutungen  $w_1$  und  $w_2$  des Pseudowortes  $w_1-w_2$  zu unterscheiden.
- Evaluere den Desambiguierer auf dem zweiten Teil durch Vergleich der Klassifikatorausgabe mit dem Originalwort.

Die	ART	die	
Wächter	NN	Wächter	
ihrerseits		ADV	ihrerseits
sind	VAFIN	sein	
nicht	PTKNEG	nicht	
hart	ADJD	hart	
,	\$,	,	
sie	PPER	sie	
haben	VAINF	haben	
ein	ART	eine	?
Herz	NN	Herz	
und	KON	und	
sind	VAFIN	sein	
nicht	PTKNEG	nicht	
ohne	APPR	ohne	
Mitgefühl		NN	Mitgefühl
.	\$.	.	

We want to extract the lemma and the POS-tag of each word.

We want a list of the lemma and a list of its tag.

Read file and generate these lists.

Die	ART	die	
Wächter	NN	Wächter	
ihrerseits		ADV	ihrerseits
sind	VAFIN	sein	
nicht	PTKNEG	nicht	
hart	ADJD	hart	
,	\$,	,	
sie	PPER	sie	
haben	VAINF	haben	
ein	ART	eine	
Herz	NN	Herz	
und	KON	und	
sind	VAFIN	sein	
nicht	PTKNEG	nicht	
ohne	APPR	ohne	
Mitgefühl		NN	Mitgefühl
.	\$.	.	

```
# Trainingsdaten einlesen
print("reading data...", file=sys.stderr)
with open(sys.argv[1]) as file:
    lemmas = []
    tags = []
    for line in file:
        word, tag, lemma = line.strip().split('\t')
        lemmas.append(lemma)
        tags.append(tag)
```

We define the pseudoword: `merged_words = ['Staat', 'Kind']`

It means we generate a new word “Statt-Kind” with `sense1 = Staat`, `sense2 = Kind`

Now we want to calculate

`f(sense1= Staat, context word)` and `f(sense2= Kind, context word)` for all context words we found in the corpus.

Die	ART	die	
Wächter	NN	Wächter	
ihrerseits		ADV	ihrerseits
sind	VAFIN	sein	
nicht	PTKNEG	nicht	
hart	ADJD	hart	
,	\$,	,	
sie	PPER	sie	
haben	VAINF	haben	
ein	ART	eine	↗
Herz	NN	Herz	
und	KON	und	
sind	VAFIN	sein	
nicht	PTKNEG	nicht	
ohne	APPR	ohne	
Mitgefühl		NN	Mitgefühl
.	\$.	.	

## How to get the context words?

- e.g. for “`sense1 = Staat`”

- iterate each word in the corpus, look for “`Staat`” (representing the word `Staat-Kind` with `sense1` ) in the corpus. The context words are 50 previous and next content words. We count their occurrences and store them in a dictionary.

Write code that computes

`count[ “Staat” ] [context] = f(“Staat”, context)`

`count [“Kind”] [context] = f(“Kind”, context)`

from lemmas and tags list



$$p(w|s) = \frac{f(s, w) + \alpha(s) p(w)}{f_1(s) + \alpha(s)}$$

Kind” oder “Staat” und  $\alpha(s)$  ist die Zahl d  
 der Nachbarschaft von  $s$  aufgetaucht sind.

$$p(w) = \frac{f_2(w)}{N}$$

Orth ufigkeit  $f_2(w)$  wie folgt definiert ist:

$$f_2(w) = f(\textit{Staat}, w) + f(\textit{Kind}, w)$$

mit  $N$  sich so ergibt:

$$N = f_1(\textit{Staat}) + f_1(\textit{Kind})$$

$$f_1(s) = \sum_w f(s, w)$$

```
# Häufigkeiten zählen
count = defaultdict(lambda: defaultdict(int))
for i in range( len(tags) ):
    # Ist das aktuelle Wort im Pseudowort enthalten?
    if lemmas[i] in merged_words:
        # Kontextwörter extrahieren
        start = max(0, i-max_dist)
        end = min(len(tags), i+max_dist+1)
        # über alle Wörter bis Maximalabstand max_dist iterieren
        for k in range(start, end):
            # falls das Kontextwort ein relevantes Tag hat und
            # nicht mit dem aktuellem Wort identisch ist
            if tags[k] in relevant_tags and k != i:
                # Häufigkeit des Paares erhöhen
                count[lemmas[i]][lemmas[k]] += 1
```

## **How to create a dictionary ?**

- `my_dict = {}`
- `my_dict = defaultdict( )`
- `my_dict = defaultdict(int)`

## **How to add a new element into the dict?**

- `my_dict[key] = value`

## **How to access the element?**

- `for key, value in my_dict.items():`  
    `print(key, value)`
- `for keys in my_dict: # or mydict.keys()`  
    `print(keys) # list of keys`
- `for values in my_dict.values():`  
    `print(values)`

## How to create a dict in dict?

- my\_dict = defaultdict(**lambda**: defaultdict(int) )
- my\_dict = defaultdict(dict )
- my\_dict = {key1: {key1\_1: value, key1\_2: value, ..},  
                  key2: {key2\_1: value, key2\_2: value,...}, ...}

## Assign new value

- my\_dict[sense][context\_word1] = 3
- my\_dict[sense][context\_word2] = 2
  
- my\_dict[sense] = {context\_word1: 3, context\_word2: 2}

## Access elements

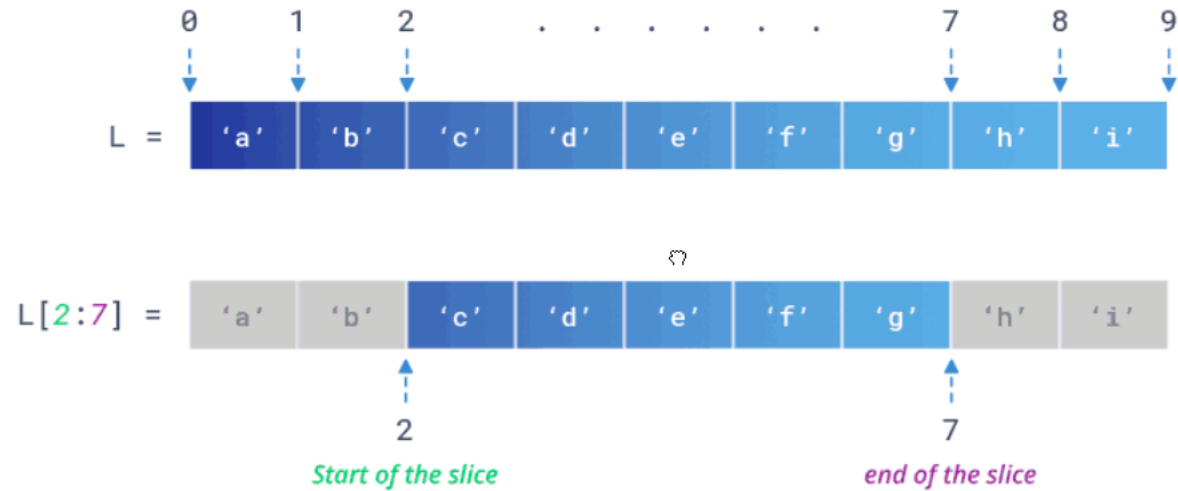
- for sense, context\_freq in my\_dict.items():  
    for context, freq in context\_freq.items():  
        print(sense, context, freq)
  
- for sense in my\_dict:  
    for context, freq in context\_freq[sense].items():  
        print(sense, context, freq)

```
from collections import defaultdict

target_dict = defaultdict(dict)
target_dict[key1][key2] = val
```

`L[start:stop:step]`

*Start position    End position    The increment*



<https://www.learnbyexample.org/python-list-slicing/>

Now we want to compute  $f(\text{sense})$  and  $f(\text{context})$

Output:  $\text{sense\_count}[\text{sense}] = f(\text{sense})$   
 $\text{cword\_count}[\text{context}] = f(\text{context})$

Using the dict "count" to compute this.

$$p(w|s) = \frac{f(s, w) + \alpha(s) p(w)}{f_1(s) + \alpha(s)}$$

Kind" oder "Staat" und  $\alpha(s)$  ist die Zahl d  
er Nachbarschaft von  $s$  aufgetaucht sind.

$$p(w) = \frac{f_2(w)}{N}$$

Worthäufigkeit  $f_2(w)$  wie folgt definiert ist:

$$f_2(w) = f(\text{Staat}, w) + f(\text{Kind}, w)$$

mit  $N$  sich so ergibt:

$$N = f_1(\text{Staat}) + f_1(\text{Kind})$$

$$f_1(s) = \sum_w f(s, w)$$

Now we want to compute  $f(\text{sense})$  and  $f(\text{context})$

Output:  $\text{sense\_count}[\text{sense}] = f(\text{sense})$   
 $\text{cword\_count}[\text{context}] = f(\text{context})$

Using the dict “count” to compute this.

```
# Gesamthäufigkeiten der Bedeutungen und Kontextwörter berechnen
sense_count = defaultdict(int)
cword_count = defaultdict(int)
for sense in count:
    for word, freq in count[sense].items():
        sense_count[sense] += freq
        cword_count[word] += freq
```

↻

$$p(w|s) = \frac{f(s, w) + \alpha(s) p(w)}{f_1(s) + \alpha(s)}$$

**Compute this prob and store it in  
logprob = { }**

Hier ist  $s$  entweder “Kind” oder “Staat” und  $\alpha(s)$  ist die Zahl der *unterschiedlichen* Inhaltswörter, die in der Nachbarschaft von  $s$  aufgetaucht sind.  $p(w)$  wird wie folgt geschätzt:

$$p(w) = \frac{f_2(w)}{N}$$

wobei die allgemeine Worthäufigkeit  $f_2(w)$  wie folgt definiert ist:

$$f_2(w) = f(\text{Staat}, w) + f(\text{Kind}, w)$$

und die Gesamthäufigkeit  $N$  sich so ergibt: ↗

$$N = f_1(\text{Staat}) + f_1(\text{Kind})$$

mit

$$f_1(s) = \sum_w f(s, w)$$



```

# Wahrscheinlichkeiten schätzen (Witten-Bell-Methode)
N = sum(sense_count.values()) # Gesamtzahl der Kontextwörter
logprob = {} # logarithmierte Wahrscheinlichkeiten der Kontextwörter
for sense in count: # für alle Lesarten
    logprob[sense] = {}
    alpha = len(count[sense]) # Zahl der Types bei den Kontextwörtern von Lesart s
    for word, freq in cword_count.items(): # für alle Kontextwörter (insgesamt)
        # Apriori-Wortwahrscheinlichkeit des Wortes berechnen
        pp = freq / N
        # mit Witten-Bell geglättete Wahrscheinlichkeit berechnen
        p = (count[sense].get(word, 0.0) + alpha * pp) / (sense_count[sense] + alpha)
        logprob[sense][word] = math.log(p)
print("done", file=sys.stderr)

# Modell in einer Datei speichern
with open(sys.argv[2], 'wt') as file:
    json.dump(logprob, file)

```

```
cat test.txt
Die      ART      die
Mutter  NN        Mutter
begleitet VVFIN    begleiten
das      ART      die
Staat:Kind NN      Staat:Kind
zur      APPRART  zu
Schule   NN        Schule
.         $.       .
```

Evaluate the model using  
this sentence

Write code that reads this  
file and generate lemmas  
and tags list

```
cat test.txt
Die      ART      die
Mutter  NN        Mutter
begleitet VVFIN   begleiten
das      ART      die
Staat:Kind NN      Staat:Kind
zur      APPRART  zu
Schule   NN        Schule
.         $.       .
```

Evaluate the model using  
this sentence

Write code that reads this  
file and generate lemmas  
and tags list

```
# zu desambiguiierende Daten einlesen
with open(sys.argv[2]) as file:
    data = [line.strip().split('\t') for line in file]
    _, tags, lemmas = zip(*data)
```

**Write code that opens the logprob from the trained model**

Write code that opens the logprob from the trained model

```
# Modell aus Datei lesen  
with open(sys.argv[1]) as file:  
    logprob = json.load( file )
```

**Write code to find the best sense for “Staat:Kind” and print the output**

Given the following variables

lemmas = list of all tokens in the input text

tags = list of the corresponding tags of each lemma

ambiguous\_word = “Staat:Kind”

senses = ['Staat', 'Kind']

max\_dist = max distance of the context words from the ambiguous word

relevant\_tags = set(['ADJA', 'ADJD', 'ADV', 'NE', 'NN', 'VVFIN', 'VVINF', 'VVPP', 'VVIZU'])

logprob[sense][word] = p(sense, context word)

```

# Vorkommen ambiger Wörter bestimmen
for i in range(len(lemmas)):
    if lemmas[i] == ambiguous_word:
        # Kontextwörter extrahieren
        start = max(0, i-max_dist)
        end = min(len(tags), i+max_dist+1)
        logp = {sense:0.0 for sense in senses}
        for k in range(start, end):
            # logarithmierte Wahrscheinlichkeiten aller
            # Kontextwörter aufsummieren
            if tags[k] in relevant_tags and k != i:
                for s in senses:
                    if lemmas[k] in logprob[s]:
                        logp[s] += logprob[s][lemmas[k]]

        # argmax-Operation berechnen
        best_sense = max(logp.keys(), key=logp.get)

        # Position und beste Lesart ausgeben
        print('Zeile %d: %s' % (i, best_sense))

```

**Aufgabe 1)** Welche Wahrscheinlichkeitsformel verwendet ein Naive-Bayes-Klassifikator? Welche Unabhängigkeitsannahmen macht der Klassifikator? Wie wählt der Klassifikator die Ergebnisklasse aus? (5 Punkte)

$$s^{\wedge} = \operatorname{argmax}_s p(s) \prod_{i=1}^n p(w_i | s)$$

Unabhängigkeitsannahme: Bag of Words

Das die Reihenfolge der Wörter egal ist  
solange sie im definierten Abstand (Fenster) auftauchen

contextwort  $w_i$  hängt nur von der Bedeutung  $s$  ab und nicht andere Kontext  
wörter

die Bedeutungsklasse mit dem höchstenwert  
Die Bedeutung mit höchsten  $p(s | C)$

**Aufgabe 2)** Wie definiert ein Naive-Bayes-Modell die gemeinsame Wahrscheinlichkeit einer Klasse (bzw. Wortlesart)  $c$  und eines Textes (bzw. Wortkontextes)  $w_1, \dots, w_n$ ?  
(2 Punkte)

Wie berechnen Sie konkret die Wahrscheinlichkeit der Klasse *Spam* für den Text: “*Diese Mail ist kein Spam !*”  
(2 Punkte)

$$p(s, C) = p(s) \prod_{i=1}^n p(w_i | s)$$

$p(c, w_1, \dots, w_n) = p(c, w_1^n) = p(c) p(w_1 | c) p(w_2 | c, w_1) p(w_3 | c, w_1, w_2) \dots$  – nicht korrekt  
 $= p(c) p(w_1 | c) p(w_2 | c) p(w_3 | c)$

Muss es bei 2) laut Formel nicht heißen:  $p(\text{Spam}, \text{Diese}, \text{Mail}, \text{ist}, \text{kein}, \text{Spam}) = p() \dots ?$   
 $p(\text{Spam}, \text{diese}, \text{Mail}, \text{ist}, \text{kein}, \text{Spam}) = p(\text{Spam}) p(\text{diese} | \text{Spam}) p(\text{Mail} | \text{Spam}) p(\text{ist} | \text{Spam}) p(\text{kein} | \text{Spam}) p(\text{Spam} | \text{Spam})$

2)  $p(\text{Spam} | \text{diese Mail ist kein Spam}) = p(\text{Spam}) p(\text{diese} | \text{Spam}) p(\text{Mail} | \text{Spam}) p(\text{ist} | \text{Spam}) p(\text{kein} | \text{Spam}) p(\text{Spam} | \text{Spam})$  --- nicht richtig



**Aufgabe 8)** Ein Naive-Bayes-Modell soll benutzt werden, um einen Zeitungsartikel einem Themengebiet (Politik, Wirtschaft etc.) zuzuweisen. Wie trainieren Sie das Modell? Welche Art von Daten benötigen Sie dazu? Wie berechnen Sie die wahrscheinlichste Kategorie eines Artikels (mit Formel)?

Anmerkung: Wir haben diese Anwendung nicht in der Vorlesung besprochen. Sie müssen hier Ihr Wissen auf eine neue Anwendung übertragen. (4 Punkte)

Wie trainieren wir das Modell?

berechnen  $\text{beste\_Thema} = \text{argmax\_themen}$  von  $p(\text{thema} | \text{Artikel})$

Daten? Korpus von Artikeln mit Themen

Wk von Kategorie eines Artikels =

$p(\text{Kategorie}, \text{Artikel}) = p(\text{Kat}) \text{ Product}_{i=1, \dots, i=n} \text{ von } p(w_i | \text{Kat})$

Artikel =  $w_1, \dots, w_n$

**Aufgabe 5)** Ein Naive-Bayes-Modell kann verwendet werden, um die wahrscheinlichste Bedeutung  $\hat{s}$  eines ambigen Wortes gegeben die Folge seiner Nachbarwörter  $w_1 \dots w_n$  zu bestimmen. Erklären Sie für jeden Schritt der folgenden Herleitung des Modelles, wie er begründet werden kann. Welche Theoreme und Definitionen werden angewandt? Welche Unabhängigkeitsannahmen werden gemacht?

$$\hat{s} = \arg \max_s p(s|w_1, \dots, w_n) \quad (1)$$

$$= \arg \max_s \frac{p(s) p(w_1, \dots, w_n|s)}{p(w_1, \dots, w_n)} \quad (2)$$

$$= \arg \max_s p(s) p(w_1, \dots, w_n|s) \quad (3)$$

$$= \arg \max_s p(s) \prod_{i=1}^n p(w_i|s, w_1, \dots, w_{i-1}) \quad (4)$$

$$= \arg \max_s p(s) \prod_{i=1}^n p(w_i|s) \quad (5)$$

$$= \arg \max_s \prod_{i=1}^n p(w_i|s) \quad (6)$$

(4 Punkte)

## Inhaltliche Wiederholung aus der Vorlesung

Naive Bayes-Modelle können zur Desambiguierung von Wortbedeutungen auf Basis der Kontextwörter verwendet werden. Die beste Bedeutung  $\hat{s}$  eines Wortes  $w$  im Kontext der Wörter  $context(w)$  wird nach der folgenden Formel berechnet:

$$\hat{s} = \arg \max_{s \in senses(w)} p(s|w) \prod_{c \in context(w)} p(c|s) \quad (1)$$

$senses(w)$ : Menge der Bedeutungen des Wort-Tokens  $w$

$context(w)$ : Liste der Kontext-Wörter

Anmerkung: Dieses Modell soll nicht nur ein ambiges Wort  $w$  sondern mehrere ambige Wörter desambiguieren.

## Aufgabe 1: Naive-Bayes-Parameterschätzung

Schreiben Sie ein Programm, welches die Wahrscheinlichkeiten  $p(s|w)$  und  $p(c|s)$  aus Trainingsdaten schätzt.

Es existiert bereits eine Funktion *extract\_data*, welche das Trainingskorpus einliest und als eine Liste von Listen zurückliefert. Ihr Programm beginnt daher mit der Zeile:

```
my @data = extract_data();
```

@data enthält eine Liste von Wortvorkommen mit der korrekten Bedeutung und den Kontextwörtern.

\$w = \$data[1][0] liefert Ihnen das 2. ambige Wort (z.B. “Bank”), das aus dem Korpus extrahiert wurde

\$s = \$data[3][1] liefert Ihnen die Bedeutung des 4. ambigen Wortes des Korpus (z.B. “Sitzbank”)

\$context = \$data[3][2] liefert Ihnen eine Referenz auf eine Liste mit den Wörtern im Kontext des 4. ambigen Wortes des Korpus (z.B. die Liste “(sitzen, Mann, vor, Baum)”)

Sie sollen die folgenden Wahrscheinlichkeiten schätzen:

- die Wahrscheinlichkeit der Bedeutung  $s$  gegeben das Wort  $w$  gemäß der Formel

$$p(s|w) = \frac{freq(w, s)}{\sum_{s'} freq(w, s')} \quad (2)$$

↗

- die Wahrscheinlichkeit des Kontextwortes  $c$  gegeben die Wortbedeutung  $s$  gemäß der Formel

$$p(c|s) = \frac{freq(s, c)}{\sum_{c'} freq(s, c')} \quad (3)$$

Schritte:

- Berechnen Sie, wie oft das ambige Wort  $w$  mit der Bedeutung  $s$  auftaucht (=  $freq(w, s)$ )
- Berechnen Sie, wie oft das Kontextwort  $c$  bei der Bedeutung  $s$  auftaucht (=  $freq(s, c)$ )
- Schätzen Sie  $p(s|w)$  nach der Formel (2)
- Schätzen Sie  $p(c|s)$  nach der Formel (3)

↗

(6 Punkte)

## Aufgabe 2: Naive-Bayes-Anwendung

Schreiben Sie nun ein Programm, welches das Naive Bayes-Modell aus einer Datei liest und auf Daten anwendet.

Ihr Programm beginnt mit

```
my(%wsprob,%scprob);  
read_params(\%wsprob,\%scprob);  
my @data = extract_data();
```

Danach sind die Parameter und die Eingabedaten eingelesen. (Die Funktionen `read_params` und `extract_data` müssen Sie nicht programmieren.)

`$wsprob{Schwester}{SCHWESTER1}` enthält die Wahrscheinlichkeit der Bedeutung „SCHWESTER1“ des Wortes „Schwester“ und entspricht  $p(s|w)$

`$scprob{SCHWESTER1}{Krankenhaus}` enthält die Wahrscheinlichkeit des Wortes „Krankenhaus“ im Kontext der Bedeutung „SCHWESTER1“ und entspricht  $\Rightarrow p(c|s)$

`@data` enthält die Eingabedaten, wobei dieses Mal die Bedeutung (z.B. `$data[3][1]`) undefiniert ist, da Sie diese ja berechnen sollen.

## 2. Wortbedeutungsdesambiguierung

Implementieren Sie einen Wortbedeutungsdesambiguierer auf Basis eines Naive-Bayes-Modelles. Auch hier schreiben Sie eine Funktion **desamb(word, tokens)**, welche das zu desambiguierende Wort und einen Text in Form einer Liste von Tokens als Argumente nimmt. Folgende Funktionen sind gegeben:

- Die Funktion **senses(word)** liefert die Liste der Bedeutungen des in “word” gespeicherten Wortes. (Bspw. könnte für das Wort *Bank* die Liste {*Bank1*, *Bank2*, *Bank3*} zurückgegeben werden.
- **wordProb(word, sense)** liefert die Wahrscheinlichkeit des Wortes “word” im Kontext der Bedeutung “sense” des zu desambiguierenden Wortes. Der Rückgabewert der Funktion ist immer größer als Null.
- **priorProb(sense)** liefert die Apriori-Wahrscheinlichkeit der Bedeutung “sense”.

Sie sollen in der Tokenliste nach Vorkommen des zu desambiguierenden Wortes suchen und diese auf Basis der Kontextwörter mit maximalem Abstand 50 desambiguieren. Ausgegeben werden soll jeweils die Position des desambiguierten Wortes und die zugewiesene Bedeutung. Sie können die Ergebnisse entweder auf die Konsole ausgeben oder in einer Liste von Paaren (Position, Bedeutung) zurückgeben.