

Hidden-Markov-Modelle: Wortart-Tagging

Hidden-Markov-Modelle: Wortart-Tagging

- Given a sequence of words, we want to find the part-of-speech (POS) of every word in the sequence.
- Some words can have more than one possible POS-tags such as “can”.
- Some words such as “I” and “a” have only one possible POS. To do POS-tagging, we can use this knowledge to narrow down the set of possible POS for these words.
- For words like “can” we still need to disambiguate their POS.
- We will use the HMM to do that.

I	can	can	a	can
<u>PRO</u>	PRO	PRO	PRO	PRO
PN	PN	PN	PN	PN
MD	MD	MD	MD	MD
NN	NN	NN	NN	NN
VB	VB	VB	VB	VB
DT	DT	DT	<u>DT</u>	DT

HMM gets a sequence of words and POS-tag sequence candidates as input and computes the best POS-tag sequence for the input.

I	want	to	shop	for	a	new	dress
PRP	VBP	TO	VB	IN	DT	JJ	NN
			NN				VB

Suppose this is our input text and the possible POS tags for each word. Here we have to disambiguate “shop” and “dress” so we will list all possible tag sequences and pass them to the HMM model.

HMM computes $p(\text{tag seq} \mid \text{word seq})$ for every tag seq candidate and does argmax to get the best tag seq, then returns it as output.

I	want	to	shop	for	a	new	dress
PRP	VBP	TO	VB	IN	DT	JJ	NN
PRP	VBP	TO	VB	IN	DT	JJ	VB
PRP	VBP	TO	NN	IN	DT	JJ	NN
PRP	VBP	TO	VB	IN	DT	JJ	VB



model



PRP VBP TO VB IN DT JJ NN

$$\hat{t}_1^n = \arg \max_{t_1^n} p(t_1^n \mid w_1^n)$$

$$\begin{aligned} p(t_1, t_2, t_3 \mid w_1, w_2, w_3) &= 0.5 \\ p(t_1, t_2, t_3 \mid w_1, w_2, w_3) &= 0.2 \\ p(t_1, t_2, t_3 \mid w_1, w_2, w_3) &= 0.2 \\ p(t_1, t_2, t_3 \mid w_1, w_2, w_3) &= 0.1 \end{aligned}$$

I	want	to	shop	for	a	new	dress
PRP	VBP	TO	VB	IN	DT	JJ	NN

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i \mid t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i \mid t_i)}_{\text{lexikalische Wk.}}$$

For every position i to $i=n+1$, we will compute the product of the context prob and the lexical prob regarding that position. Then, multiply every result together to get the final probability.

The model estimates $p(t_1, t_2, t_3 \mid w_1, w_2, w_3)$, based on two things.

- **Lexical probability**
 - how likely it is that “I” would have the tag “PRP” ?
 - how likely it is that “shop” would have the tag “VB” ?
- **Context probability** (Let’s say here we consider the previous 2 tags)
 - if the previous 2 tags are “VBP, TO” what is the probability that VB will be the tag at position 4(shop)? (is it more than having NN?)
- These probabilities (=model parameters) are estimated by looking in a corpus.
- If these probabilities are high, the total probability will be high.

I	want	to	shop	for	a	new	dress
PRP	VBP	TO	VB	IN	DT	JJ	NN
1.0	1.0	1.0	0.6	1.0	1.0	1.0	0.5

This is just the theoretical concept. In practice, we do not compute the best tag seq this way because it is inefficient --> we will use the Viterbi algorithm instead

$$\hat{t}_1^n = \arg \max_{t_1^n} p(t_1^n | w_1^n)$$

$$\arg \max_{t_1^n} \frac{p(t_1^n, w_1^n)}{p(w_1^n)}$$

bedingte Wahrscheinlichkeit: $p(x|y) = \frac{p(x,y)}{p(y)}$

$$\arg \max_{t_1^n} p(t_1^n, w_1^n)$$

Die Konstante $p(w_1^n)$ hat keinen Einfluss auf das Maximierungsergebnis.

$$p(t_1^n, w_1^n) = \left(\prod_{i=1}^{n+1} p(t_i | t_1^{i-1}) \right) \prod_{i=1}^{n+1} p(w_i | t_1^{n+1}, w_1^{i-1})$$

Zerlegung in ein Produkt bedingter Wahrscheinlichkeiten:

Kettenregel $p(x, y, z) = p(x)p(y|x)p(z|xy)$

Example

Wir fügen hier ein Endetag $t_{n+1} = \langle /s \rangle$ und ein Endetoken $w_{n+1} = \epsilon$ hinzu,
mit $p(\epsilon | \langle /s \rangle) = 1$.

$$p(t_1, t_2, t_3, w_1, w_2, w_3) = p(t_1) p(t_2 | t_1) p(t_3 | t_1, t_2) p(\langle /s \rangle | t_1, t_2, t_3) \\ p(w_1 | t_1, t_2, t_3, \langle /s \rangle) p(w_2 | t_1, t_2, t_3, \langle /s \rangle, w_1) p(w_3 | t_1, t_2, t_3, \langle /s \rangle, w_1, w_2) p(\text{leer} | t_1, t_2, t_3, \langle /s \rangle, w_1, w_2, w_3)$$

w1	w2	w3	e
t1	t2	t3	</s>

Model Herleitung (2/2)

$$p(t_1^n, w_1^n) = \left(\prod_{i=1}^{n+1} p(t_i | t_1^{i-1}) \right) \prod_{i=1}^{n+1} p(w_i | t_1^{n+1}, w_1^{i-1})$$

move the product symbol to the front

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} p(t_i | t_1^{i-1}) p(w_i | w_1^{i-1}, t_1^{n+1})$$

just switch the position

Wir machen nun die folgenden vereinfachenden Annahmen:

- Das Wortart-Tag t_i hängt nur von den k vorherigen Tags ab.
- Das Wort w_i hängt nur von seiner Wortart t_i ab.
- Die Wahrscheinlichkeiten sind unabhängig von der Position.

the final formula

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

Dieses Modell heißt **Hidden-Markow-Modell**, weil die Zustände (Tags bzw. Tagpaare beim Trigramm-Tagger) nicht direkt beobachtbar sind.

Wir fügen am Anfang k Starttags hinzu, damit $p(t_1 | t_{1-k}^0)$ definiert ist.

$$p(t_1, t_2, t_3, w_1, w_2, w_3) = p(t_1) p(t_2 | t_1) p(t_3 | t_1, t_2) p(</s> | t_1, t_2, t_3) \\ p(w_1 | t_1, t_2, t_3, </s>) p(w_2 | t_1, t_2, t_3, </s>, w_1) \\ p(w_3 | t_1, t_2, t_3, </s>, w_1, w_2) p(\text{leer} | t_1, t_2, t_3, </s>, w_1, w_2, w_3)$$

$$p(t_1, t_2, t_3, w_1, w_2, w_3) = p(t_1 | <s>, <s>) p(t_2 | <s>, t_1) p(t_3 | t_1, t_2) p(</s> | t_2, t_3) \\ p(w_1 | t_1, t_2, t_3, </s>) p(w_2 | t_1, t_2, t_3, </s>, w_1) \\ p(w_3 | t_1, t_2, t_3, </s>, w_1, w_2) \\ p(\text{leer} | t_1, t_2, t_3, </s>, w_1, w_2, w_3)$$

k=2

e	e	w1	w2	w3	e
<s>	<s>	t1	t2	t3	</s>

$$\hat{t}_1^n = \arg \max_{t_1^n} p(t_1^n | w_1^n)$$

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

Zwei Typen von Parametern

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

Die **Kontextwahrscheinlichkeiten** erfassen die syntaktischen Abhängigkeiten der Wortart-Tags vom linken Satzkontext.

Die **lexikalischen Wahrscheinlichkeiten** erfassen die syntaktischen Eigenschaften des aktuellen Wortes.

Beispiel

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} p(t_i | t_{i-k}^{i-1}) p(w_i | t_i)$$

Bei einem HMM 1. Ordnung (k=1) gilt:

$$p(\text{NE VVFIN, Peter liest}) = p(\text{NE} | \langle s \rangle) p(\text{Peter} | \text{NE}) \cdot \\ p(\text{VVFIN} | \text{NE}) p(\text{liest} | \text{VVFIN}) \cdot \\ p(\langle /s \rangle | \text{VVFIN}) p(\epsilon | \langle /s \rangle)$$

Beim HMM 2. Ordnung (k=2) gilt:

$$p(\text{NE VVFIN, Peter liest}) = p(\text{NE} | \langle s \rangle, \langle s \rangle) p(\text{Peter} | \text{NE}) \cdot \\ p(\text{VVFIN} | \langle s \rangle, \text{NE}) p(\text{liest} | \text{VVFIN}) \cdot \\ p(\langle /s \rangle | \text{NE}, \text{VVFIN}) p(\epsilon | \langle /s \rangle)$$

pos 0	1	2	3
e	Peter	liest	e
<s>	NE	VVFIN	</s>

pos1:

-lexical: $p(\text{Peter} | \text{NE})$
-context: $p(\text{NE} | \langle s \rangle)$

pos2

-lexical: $p(\text{liest} | \text{VVFIN})$
-context: $p(\text{VVFIN} | \text{NE})$

pos3

-lexical: $p(\epsilon | \langle /s \rangle)$
-context: $p(\langle /s \rangle | \text{VVFIN})$

Final answer

$$p(\text{Peter} | \text{NE}) p(\text{NE} | \langle s \rangle) \cdot \\ p(\text{liest} | \text{VVFIN}) p(\text{VVFIN} | \text{NE}) \cdot \\ p(\epsilon | \langle /s \rangle) p(\langle /s \rangle | \text{VVFIN})$$

1

2

3

Aufgabe 3) Wie berechnen Sie bei einem **Hidden-Markowmodell** **2. Ordnung** die gemeinsame Wahrscheinlichkeit der Wortfolge *Bayern, schlägt, Manchester* und der Tagfolge *NE, VVFİN, NE*? (2 Punkte)

K=2

e e Bayern, schlägt, Manchester e
 <s> <s> NN, VVFİN, NE </s>

$$p(\text{Bayern, schlägt, Manchester, NN, VVFİN, NE}) = p(\text{NN} \mid \langle s \rangle, \langle s \rangle) p(\text{Bayern} \mid \text{NN}) \\
p(\text{VVFİN} \mid \langle s \rangle, \text{NN}) p(\text{schlägt} \mid \text{VVFİN}) \\
p(\text{NE} \mid \text{NN, VVFİN}) p(\text{Manchester} \mid \text{NE}) \\
p(\langle /s \rangle \mid \text{VVFİN, NE}) p(e \mid \langle /s \rangle)$$

context: $p(\text{NN} \mid \langle s \rangle, \langle s \rangle)$

lexical: $p(\text{Bayern} \mid \text{NN})$

Aufgabe 2) Wie berechnen Sie bei einem **Markowmodell** 2. Ordnung (Trigramm-Modell) über Buchstaben die Wahrscheinlichkeit der Buchstabenfolge *CIS*? (2 Punkte)

$k=2$

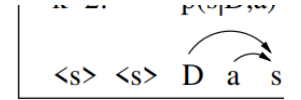
$\langle s \rangle \langle s \rangle C I S \langle /s \rangle$

$$p(C, I, S) = p(C | \langle s \rangle, \langle s \rangle) * p(I | \langle s \rangle, C) * p(S | C, I) p(\langle /s \rangle | I, S)$$

example

note: in the character-level model, **spaces** should also be seen as a symbol.

$$p_L(T) = \prod_{i=1}^{n+1} p_L(a_i | a_{i-k} \dots a_{i-1})$$



Ein solches Modell heißt **Markowmodell** der Ordnung k .

Für $k = 2$ bekommen wir:

$$p(\text{Er ölt}) = p(E | \langle s \rangle, \langle s \rangle) p(r | \langle s \rangle, E) p(- | E, r) p(\text{ö} | r, -) p(I | -, \text{ö}) p(t | \text{ö}, I) p(\langle /s \rangle | I, t)$$

⇒ Es werden k Startsymbole $\langle s \rangle$ und ein Endesymbol $\langle /s \rangle$ hinzugefügt:
 $a_{-1} = a_0 = \langle s \rangle$ und $a_{n+1} = \langle /s \rangle$

Bei einem HMM 1. Ordnung (k=1) gilt:

$$p(\text{NE VVFIN, Peter liest}) = p(\text{NE}|\langle s \rangle) p(\text{Peter}|\text{NE}) \cdot p(\text{VVFIN}|\text{NE}) p(\text{liest}|\text{VVFIN}) \cdot p(\langle /s \rangle|\text{VVFIN}) p(\epsilon|\langle /s \rangle)$$

Parameterschätzung

$$p(t|t_1, \dots, t_k) : \quad p(w|t)$$

Zur Schätzung der Parameter braucht man ein **Trainingskorpus**, in dem jedes Wort mit seiner Wortart annotiert ist.

Man zählt die k+1-Gramm-Häufigkeiten $f(t_1, \dots, t_{k+1})$ und die Tag-Wort-Häufigkeiten $f(t, w)$ und schätzt dann die Wahrscheinlichkeiten wie folgt:

$$p(w|t) = \frac{f(t, w)}{\sum_{w'} f(t, w')} \quad p(t|t_1, \dots, t_k) = \frac{f(t_1, \dots, t_k, t)}{\sum_{t'} f(t_1, \dots, t_k, t')}$$

Beispiel:

$$p(\text{Haus}|\text{NN}) = \frac{f(\text{NN, Haus})}{\sum_{w'} f(\text{NN, } w')} \quad p(\text{NN}|\text{ART}) = \frac{f(\text{ART, NN})}{\sum_{t'} f(\text{ART, } t')}$$

Die Kontextwahrscheinlichkeiten $p(t|t_1, \dots, t_k)$ müssen aber noch geglättet werden.

Training data (each word is annotated with its POS)

Ich wohne in meinem Haus.

POS, ... ART NN

Er hat ein Haus gekauft.

... ART NN.....

Ein Hund ist gut.

ART NN V ADJ

Extract the frequency of (tag, word) and tag n-gram (tag1, tag2)



f(Tag, word)

NN, Haus: 100
NN, Hund: 150



f(Tag1, Tag2)

NN, V : 9
NN, ART : 3
ART, NN : 8



Estimate the probabilities

$p(w|t)$

$p(t|t_1, \dots, t_k) :$

$$\hat{t}_1^n = \arg \max_{t_1^n} p(t_1^n | w_1^n)$$

Formulate the problem (compute the best tag seq between tag sequence candidates)

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

- Define how to compute the prob $p(\text{tag_seq} | \text{word_seq})$
- Here we choose HMM

$$p(w|t) = \frac{f(t, w)}{\sum_{w'} f(t, w')} \quad p(t|t_1, \dots, t_k) = \frac{f(t_1, \dots, t_k, t)}{\sum_{t'} f(t_1, \dots, t_k, t')}$$

- Choose a method to estimate the probs
- Here we choose the relative frequency method

Better : apply a smoothing method such as backoff smoothing

Model training

After we defined how p should be computed, we can collect relevant information to compute them. E.g. choose $k=1$, then we extract tag bigrams.

Training data (each word is annotated with its POS)
Ich wohne in meinem Haus.
POS, ... ART NN
Er hat ein Haus gekauft.
... ART NN.....
Ein Hund ist gut.
ART NN V ADJ

Extract the frequency of (tag, word) and tag n-gram (tag1, tag2)



$f(\text{Tag}, \text{word})$

NN, Haus: 100
NN,Hund: 150



$f(\text{Tag1}, \text{Tag2})$

NN, V : 9
NN,ART : 3
ART, NN : 8



Estimate the probabilities

$p(w|t)$

$p(t|t_1, \dots, t_k)$

Prediction

Use the model to predict the tag sequence of a new text.

„Peter, liest“

list all possible tags for „Peter“ and „liest“

compute $p(\text{tags} | \text{text})$ for every possible tag seq,
return the best tag seq

Search for the best tag seq

Suche

Problem: Es gibt zuviele mögliche Tagfolgen für eine gegebene Wortfolge, um alle aufzählen und ihre Wahrscheinlichkeit berechnen zu können.

Beobachtung: Wenn zwei mögliche Wortartfolgen $T = t_1, \dots, t_m$ und $S = s_1, \dots, s_m$ (für die ersten m Wörter eines Satzes) in den letzten k Tags übereinstimmen, dann kann bei einem HMM k -ter Ordnung die weniger wahrscheinliche der beiden Tagfolgen nicht ein Präfix der besten Gesamtagfolge sein, und wir können sie ignorieren.

Widerspruchsbeweis

Angenommen S ist weniger wahrscheinlich als T , aber ein Präfix der besten Tagfolge. Dann hat die beste Tagfolge die Form SV , wobei V eine weitere Tagfolge ist. In diesem Fall hätte aber die Tagfolge TV eine höhere Wahrscheinlichkeit als SV , weil T wahrscheinlicher als S ist und die Wahrscheinlichkeiten, die jeweils für V hinzumultipliziert werden, bei SV und TV identisch sind. Also kann S kein Präfix der besten Tagfolge sein. \square

Beispiel: $k=2$ $T=X,A,B$ $S=Y,A,B$ $V=C,D$
 $TV=X,A,B,C,D$ $SV=Y,A,B,C,D$

C hängt hier nur von A und B ab. D hängt nur von B und C ab.
Also sind die Tags C,D nach X,A,B genauso wahrscheinlich wie nach Y,A,B .

Suche

Diese Eigenschaft erlaubt eine effiziente Verarbeitung, hat aber zur Folge, dass manche Ambiguitäten nicht korrekt aufgelöst werden können:

The	horse	raced	past	the	barn	(fell)	.
DT	NN	VBD	IN	DT	NN		
DT	NN	VCN	IN	DT	NN		

Ein HMM 2. Ordnung würde das Wort "raced" gleich taggen, egal ob "fell" nachfolgt oder nicht. Erst ein HMM 4. Ordnung kann weit genug zurückschauen, um korrekt zu desambiguieren.

Anmerkung: Man könnte das Problem lösen, indem man das Tagset verfeinert und bspw. alle Tags, die nach einem finiten Verb folgen, mit einem Apostroph markiert.

The	horse	raced	past	the	barn	(fell)	.
DT	NN	VBD	IN'	DT'	NN'		\$.
DT	NN	VCN	IN	DT	NN	VBD	\$.

$$\hat{t}_1^n = \arg \max_{t_1^n} p(t_1^n | w_1^n)$$

The initial way to determine the best tag seq is not optimal because

1. there can be too many tag seq candidates (require a lot of computing)
2. The way HMM computes $p(\text{tags} | \text{words})$ is not good because it can assign low prob to good tag seq. This can happen when, e.g., two candidates have the same tag suffix.

Suppose, Y,A,B,C,D is the best tag. Another tag seq X,A,B,C,D can have higher prob than the best tag if its prefix X,A,B has a higher prob than Y,A,B

Solution: Viterbi

Viterbi-Algorithmus (Bigramm-Tagger)

1. Initialisierung: $\delta_t(0) = \begin{cases} 1 & \text{falls } t = \langle s \rangle \\ 0 & \text{sonst} \end{cases}$

2. Berechnung: (für $0 < k \leq n + 1$)

$$\delta_t(k) = \max_{t'} \delta_{t'}(k-1) p(t|t') p(w_k|t)$$

$$\psi_t(k) = \arg \max_{t'} \delta_{t'}(k-1) p(t|t') p(w_k|t)$$

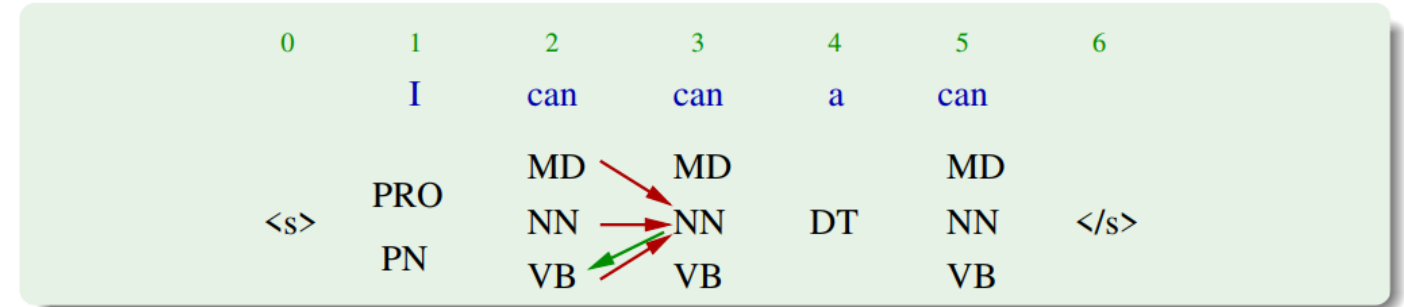
3. Ausgabe: (für $0 < k \leq n$)

$$t_{n+1} = \langle /s \rangle$$

$$t_k = \psi_{t_{k+1}}(k+1)$$

t, t' sind Tags, k sind Wortpositionen, $\delta_t(k)$ sind die Viterbiwahrscheinlichkeiten und $\psi_t(k)$ sind die besten Vorgängertags

Beispiel: Viterbi-Algorithmus



$$\delta_{NN}(3) = \max \begin{pmatrix} \delta_{MD}(2) p(NN|MD) p(can|NN), \\ \delta_{NN}(2) p(NN|NN) p(can|NN), \\ \delta_{VB}(2) p(NN|VB) p(can|NN) \end{pmatrix}$$

(I will talk about it in the next session)

Smoothing

$$\hat{t}_1^n = \arg \max_{t_1^n} p(t_1^n | w_1^n)$$

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

$$p(w|t) = \frac{f(t, w)}{\sum_{w'} f(t, w')}$$

Lexikal. WK

$$p(t|t_1, \dots, t_k) = \frac{f(t_1, \dots, t_k, t)}{\sum_{t'} f(t_1, \dots, t_k, t')}$$

Kontext WK

- Choose a method to estimate the probs
- Here we choose the relative frequency method

Better : apply a smoothing method such as backoff smoothing

Now, we want to smooth these probabilities

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

1) use the backoff smoothing to smooth the context prob

See the detail in Übung7

- Kontextwahrscheinlichkeiten $p(t|t', t'')$

Diese Wahrscheinlichkeiten sollen mit Kneser-Ney-Glättung geglättet werden. Sie müssen hier nur die relativen Häufigkeiten $p^*(t|t', t'')$, $p^*(t|t'')$ und $p^*(t)$ (vgl. Übung 3) berechnen. $p(t|t', t'')$ wird erst im eigentlichen Taggerprogramm berechnet.

$$p(t_1^n, w_1^n) = \prod_{i=1}^{n+1} \underbrace{p(t_i | t_{i-k}^{i-1})}_{\text{Kontextwahrsch.}} \underbrace{p(w_i | t_i)}_{\text{lexikalische Wk.}}$$

2) smooth the lex. prob

adjust the formula

$$p(w|t) = \frac{p(t|w)p(w)}{p(t)}$$

$p(w)$ is constant, so it gets canceled out.

Die Apriori-Wahrscheinlichkeit $p(t) = f(t)/N$ (N = Korpusgröße) der Tags kann leicht aus dem Trainingskorpus geschätzt werden.

backoff $p(t|w)$ by using the suffix and the case (Groß-/Kleinschreibung) of the word.

für Suffixlänge $l=5$:

$$p(\text{NN}|\text{schwärzlich}) = r(\text{NN}|\text{schwärzlich}) + \alpha(\text{schwärzlich}) ($$

$$\frac{r_{\text{suffix}}(\text{NN}|\text{zlichK}) + \alpha_{\text{suffix}}(\text{zlichK}) ($$

$$r_{\text{suffix}}(\text{NN}|\text{lichK}) + \alpha_{\text{suffix}}(\text{lichK}) ($$

$$r_{\text{suffix}}(\text{NN}|\text{ichK}) + \alpha_{\text{suffix}}(\text{ichK}) ($$

$$r_{\text{suffix}}(\text{NN}|\text{chK}) + \alpha_{\text{suffix}}(\text{chK}) ($$

$$r_{\text{suffix}}(\text{NN}|\text{hK}) + \alpha_{\text{suffix}}(\text{hK}) p(\text{NN}|\text{K}))))))$$

the recursion ends here

This helps prevent zero prob when the new text contains unseen words or tags

Behandlung unbekannter Wörter

Die Apriori-Wahrscheinlichkeit $p(t) = f(t)/N$ (N = Korpusgröße) der Tags kann leicht aus dem Trainingskorpus geschätzt werden.

Die bedingte Wahrscheinlichkeit $p(t|w)$ kann anhand des Wortsuffixes (z.B. der Länge l) geschätzt werden:

$$p(t|w) = p(t|a_1 a_2 \dots a_n) \approx p_{\text{suff}}(t|a_{n-l+1} \dots a_n)$$

$$p(\text{NN}|\text{schwärzlich}) \approx p_{\text{suff}}(\text{NN}|\text{zlich}) \quad \text{falls } l = 5$$

Die Suffix-Tag-Wahrscheinlichkeit (α ist hier eine beliebige Buchstabenfolge.)

$$p_{\text{suff}}(t|s_1 \dots s_l) = \frac{f_{\text{suff}}(s_1 \dots s_l, t)}{\sum_{t'} f_{\text{suff}}(s_1 \dots s_l, t')} \quad f_{\text{suff}}(s_1 \dots s_l, t) = \sum_{w=\alpha s_1 \dots s_l} f(w, t)$$

wird aus Trainingsdaten geschätzt und mit einem Backoff-Verfahren geglättet.

Für groß- und kleingeschriebene Wörter schätzt man am besten separate Parameter. Dafür kann man einfach einen Buchstaben (z.B. "G" für Groß- und "K" für Kleinschreibung) an das Wortsuffix anhängen: $p_{\text{suff}}(\text{NN}|\text{zlichK})$

Navigationssymbole

Behandlung unbekannter Tags

neue Formel zum Taggen mit unbekannten Wörtern:

$$\hat{t}_1^n = \arg \max_{t_1^n} \prod_{i=1}^{n+1} p(t_i | t_{i-k}^{i-1}) \frac{p(t_i | w_i)}{p(t_i)}$$

Neben unbekannten Wörtern stellen auch **ungesehene Tags** bekannter Wörter ein Problem dar (bspw. wenn das englische Wort "indicate" nur als Infinitiv aber nicht als finites Verb aufgetaucht ist.)

Die suffixbasierten Tag-Wahrscheinlichkeiten geben meist auch den ungesehenen Tags eine kleine Wahrscheinlichkeit.

Daher ist es sinnvoll, wortbasierte und suffixbasierte Wahrscheinlichkeiten in einem Backoff-Modell zu kombinieren:

$$p(t|w) = r(t|w) + \alpha(w) p_{\text{suff}}(t|\text{suffix}(w))$$

$$\text{mit } r(t|w) = \max(0, \frac{f(w, t) - \delta}{\sum_{t'} f(w, t')})$$

Navigationssymbole

Hidden Markov Model (HMM) is a [statistical Markov model](#) in which the system being [modeled](#) is assumed to be a [Markov process](#) — call it X — with unobservable ("*hidden*") states. As part of the definition, HMM requires that there be an observable process Y whose outcomes are "influenced" by the outcomes of X in a known way. Since X cannot be observed directly, the goal is to learn about X by observing Y . HMM has an additional requirement that the outcome of Y at time $t = t_0$ may be "influenced" exclusively by the outcome of X at $t = t_0$ and that the outcomes of X and Y at $t < t_0$ must not affect the outcome of Y at $t = t_0$.

Hidden Markov models are known for their applications to [thermodynamics](#), [statistical mechanics](#), [physics](#), [chemistry](#), [economics](#), [finance](#), [signal processing](#), [information theory](#), [pattern recognition](#) - such as [speech](#), [handwriting](#), [gesture recognition](#),^[1] [part-of-speech tagging](#), [musical score following](#),^[2] [partial discharges](#)^[3] and [bioinformatics](#).^{[4][5]}