



Computational Linguistics 2014-2015

- **Walter Daelemans** (walter.daelemans@uantwerpen.be)
- **Guy De Pauw** (guy.depauw@uantwerpen.be)
- **Mike Kestemont** (mike.kestemont@uantwerpen.be)

<http://www.clips.uantwerpen.be/cl1415>



Practical

Location	P0.11 (Scribanihuis)
Reading material	<ul style="list-style-type: none">• D. Jurafsky & J.H. Martin (2009) <i>Speech and Language Processing - An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition</i> (2nd ed). Pearson Education, USA.• Natural Language Processing with Python
Software	Python 3.4 and NLTK: Installation Instructions
Evaluation	Take-home assignments and oral examination
Lecturers	Walter Daelemans: walter.daelemans@uantwerpen.be Mike Kestemont: mike.kestemont@uantwerpen.be Guy De Pauw: guy.depauw@uantwerpen.be



Program

Session	Day	Date	Chapter	Topic	Reading Assignment	Slides	Take-home Assignment	
1	Monday	29/9/2014	Python	Session 1 - Variables	See Github			
2	Thursday	2/10/2014	Python	Session 2 - Collections				
3	Monday	6/10/2014	Python	Session 3 - Conditions (and an introduction to loops)				
4	Thursday	9/10/2014	Python	Session 4 - Loops				
5	Monday	13/10/2014	Python	Session 5 - Reading and writing to files				
6	Thursday	16/10/2014	Python	Session 6 - Writing your own Functions and importing packages				
7	Monday	20/10/2014	Python	Session 7 - Regular Expressions in Python				
8	Thursday	23/10/2014	Python	Session 8 - Advanced looping in Python and list comprehensions				
9	Monday	27/10/2014	Theory	Introduction to Computational Linguistics	Jurafsky & Martin: Chapter 1	PDF		
10	Monday	3/11/2014	Theory	Regular Expressions and Finite State Automata & Transducers	Jurafsky & Martin: Chapter 2 ; Chapter 3	Slides morfsegment.py	See last slide. Deadline: 24/11 participles.py	
	Monday	10/11/2014	Remembrance day: no session					
11	Monday	17/11/2014	Theory	Part-of-Speech Tagging	Jurafsky & Martin: Chapter 5 (not 5.5, 5.8 and 5.9)	Slides Python Code	See last slide. Deadline: 8/12	
12	Monday	24/11/2014	Theory	Syntactic Analysis & Parsing	Jurafsky & Martin: Chapter 12 (not 12.7.2, 12.8); Chapter 13 (not 13.4.1, 13.4.2, 13.5.1)			
13	Monday	1/12/2014	Theory	Minimum Edit Distance + Probabilistic Methods	Jurafsky & Martin: Chapter 3.11; Chapter 4.1, 4.2 and 4.3; Chapter 5.5 and 5.9; Chapter 14.1, 14.3 and 14.4;			
14	Monday	8/12/2014	Theory	Word Sense Disambiguation	Jurafsky & Martin: Chapter 19.1, 19.2, 19.3, Chapter 20 (20.1->20.5)			
15	Monday	15/12/2014	Theory	Sentence semantics and discourse; Information extraction	Jurafsky & Martin: Chapter 21; Chapter 22			



Bayesian Inference

N-gram models



Statistical Methods

- automatically derive statistical data from (annotated) corpora
- frequency of observed events are interpreted as the **probability** of those events occurring in the future
- We can use these probabilities to perform disambiguation

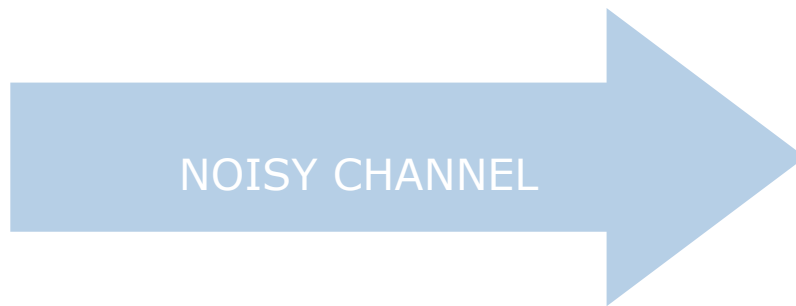
e.g. Most likely tag for *can* in "I can do this."?

$P(\text{MD}|\text{can})$ vs $P(\text{NN}|\text{can})$ vs $P(\text{VB}|\text{can})$

- $P(x|y)$ is calculated through *Bayesian Inference*



INPUT



Noisy Channel

OUTPUT

<u>Task</u>	<u>Input</u>	<u>Output</u>
Speech Recognition	String of Words	Acoustic Signal
OCR/ Spellchecking	Correct Text	Text with errors
POS Tagging	String of POS Tags	String of words
Machine Translation	Sentence in English	Sentence in Chinese

$P(\text{Input}|\text{Output})$?????



Bayesian Inference

Bayes' Law

$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$

e.g. From Wikipedia
Drug Test: 0.99 accurate (99% chance that a user tests positive, 99% chance that a non-user tests negative)
Users: 0.5% of the population

What is the probability that someone who tests positive, is a user?

$$\begin{aligned} P(\text{User}|+) &= \frac{P(+|\text{User}) \cdot P(\text{user})}{P(+)} \\ &= \frac{0.99 * 0.005}{P(+|\text{User}) * P(\text{User}) + P(+|\text{non-user}) * P(\text{non-user})} \\ &= \frac{0.99 * 0.005}{0.99 * 0.005 + 0.01 * 0.995} \\ &= 0.332 \end{aligned}$$



Bayesian Inference

Bayes' Law

$$P(x|y) = \frac{P(y|x).P(x)}{P(y)}$$

From a corpus we calculated the following probabilities

$P(\text{can}|\text{MD}) = 0.8$ (the frequency with which *can* was observed as MD)

$P(\text{can}|\text{NN}) = 0.1$ and $P(\text{can}|\text{VB}) = 0.1$

$P(\text{MD}) = 0.05$, $P(\text{NN}) = 0.3$ and $P(\text{VB}) = 0.1$

$P(\text{can}) = 0.00001$

What is $P(\text{MD}|\text{can})$, the probability that we need to tag MD when we see 'can'?

$$\begin{aligned} P(\text{MD}|\text{can}) &= \frac{P(\text{can}|\text{MD}).P(\text{MD})}{P(\text{can})} \\ &= \frac{0.8 * 0.05}{1} = 0.04 \end{aligned}$$

$$P(\text{NN}|\text{can}) = P(\text{can}|\text{NN}).P(\text{NN}) = 0.1 * 0.3 = 0.03$$

$$P(\text{VB}|\text{can}) = P(\text{can}|\text{VB}).P(\text{VB}) = 0.1 * 0.1 = 0.01$$



$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$

"can" is counted 60 times as "MD" in corpusA and 40 times as "NN". In corpusB "can" is counted 70 times as "NN" and 30 times as "MD".

1. What is the probability of "can" as "MD" in corpusA?
2. What is the probability of "can" as "NN" in corpusB?
3. We pick a sentence randomly from one of the 2 corpora:
"I can do this"

What is the probability that this sentence came from corpusA?



$$P(x|y) = \frac{P(y|x) \cdot P(x)}{P(y)}$$

"can" is counted 60 times as "MD" in corpusA and 40 times as "NN". In corpusB "can" is counted 70 times as "NN" and 30 times as "MD".

1. What is the probability of "can" as "MD" in corpusA?
2. What is the probability of "can" as "NN" in corpusB?
3. We pick a sentence randomly from one of the 2 corpora:
"I can do this"

What is the probability that this sentence came from corpusA?

$$\begin{aligned} P(\text{corpusA}|\text{canMD}) &= (P(\text{canMD}|\text{corpusA}) \cdot P(\text{corpusA})) / P(\text{canMD}) \\ &= (60/100 \times 1/2) / (60+30/200) \\ &= (0.6 \times 0.5) / 0.45 \\ &= 0.667 \end{aligned}$$



Bayesian Inference

In language technology, we calculate the probability of the association between an input sequence and an output sequence.

e.g. Machine translation

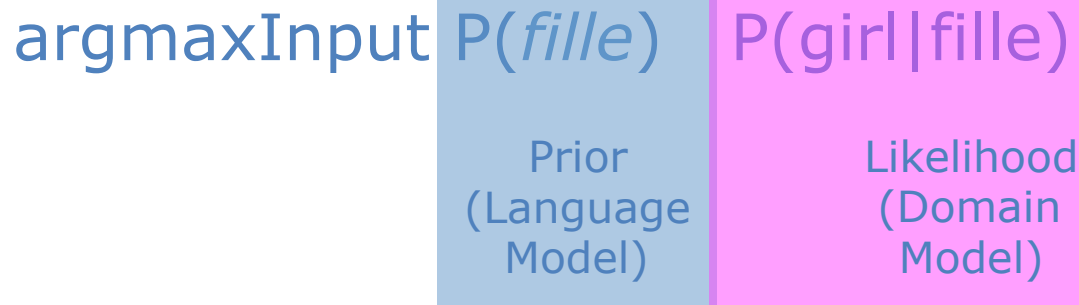
$\text{argmax}_{\text{Input}} P(\textit{fille} | \textit{girl}) =$

$\text{argmax}_{\text{Input}}$	$P(\textit{fille})$	$P(\textit{girl} \textit{fille})$
	Prior (Language Model)	Likelihood (Domain Model)



Bayesian Inference

$\text{argmaxInput } P(\textit{fille}|\textit{girl}) =$



The Domain Model provides the probability that *girl* can be translated as *fille*

The language model provides the probability that the word *fille* exist (in that context)



Bayes' Rule & Noisy Channel

	$P(\text{Input})$	$P(\text{Output} \text{Input})$
Machine Translation	Language Model	Translation model
OCR		Model of OCR errors
Spellchecking		Model of spelling errors
POS-Tagging		Tag-Word Model
Speech Recognition		Acoustic model



- What is the probability of a given sequence of words, tokens, tags?
- Most common: n-gram models
- data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$
- unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
 $P(\text{sentence}) = \prod P(\text{word})$



- What is the probability of a given sequence of words, tokens, tags?
- Most common: n-gram models
- data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$
- unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
 $P(\text{sentence}) = \prod P(\text{word})$



- What is the probability of a given sequence of words, tokens, tags?
- Most common: n-gram models
- data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$
- unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
 $P(\text{sentence}) = \prod P(\text{word})$



- What is the probability of a given sequence of words, tokens, tags?
- Most common: n-gram models
- data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$
- unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
 $P(\text{sentence}) = \prod P(\text{word})$



- But unigram is a weak language model
- Suppose we want to predict the most likely possible word in the sentence

Just then, the white...

According to unigram:

$$P(\textit{the}) = 0.07$$

$$P(\textit{rabbit}) = 0.00001$$

And so

$$P(\textit{Just then, the white the}) > P(\textit{Just then, the white rabbit})$$

Although intuitively

$$P(\textit{Just then, the white the}) < P(\textit{Just then, the white rabbit})$$

- Contextual information limited to n-value (cfr. n-gram models)



Language Model

- $P(\text{sentence}) = \prod P(\text{word})$
- Unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
- bigram: $P(\text{word}_i | \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-1})$

$P(\text{rabbit} | \text{white}) = \text{freq}(\text{white rabbit}) / \text{freq}(\text{white})$

$P(\text{the} | \text{white}) = \text{freq}(\text{white the}) / \text{freq}(\text{the})$

data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$



Language Model

- $P(\text{sentence}) = \prod P(\text{word})$
- Unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
- bigram: $P(\text{word}_i | \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-1})$

$P(\text{rabbit} | \text{white}) = \text{freq}(\text{white rabbit}) / \text{freq}(\text{white})$

$P(\text{the} | \text{white}) = \text{freq}(\text{white the}) / \text{freq}(\text{the})$

data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$



Language Model

- $P(\text{sentence}) = \prod P(\text{word})$
- Unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
- bigram: $P(\text{word}_i | \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-1})$

$$P(\text{rabbit} | \text{white}) = \text{freq}(\text{white rabbit}) / \text{freq}(\text{white})$$

$$P(\text{the} | \text{white}) = \text{freq}(\text{white the}) / \text{freq}(\text{the})$$

data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$



Language Model

- $P(\text{sentence}) = \prod P(\text{word})$
- Unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
- bigram: $P(\text{word}_i | \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-1})$
- trigram: $P(\text{word}_i | \text{word}_{i-2} \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-2} \text{word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-2} \text{word}_{i-1})$

$P(\text{rabbit} | \text{the white}) = \text{freq}(\text{the white rabbit}) / \text{freq}(\text{the white})$

$P(\text{the} | \text{the white}) = \text{freq}(\text{the white the}) / \text{freq}(\text{the white})$

data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$



- $P(\text{sentence}) = \prod P(\text{word})$
- Unigram: $P(\text{word}) = \text{freq}(\text{word}) / N$
- bigram: $P(\text{word}_i | \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-1})$
- trigram: $P(\text{word}_i | \text{word}_{i-2} \text{word}_{i-1}) = \text{freq}(\text{'word}_{i-2} \text{word}_{i-1} \text{word}_i\text{'}) / \text{freq}(\text{word}_{i-2} \text{word}_{i-1})$

$P(\text{rabbit} | \text{the white}) = \text{freq}(\text{the white rabbit}) / \text{freq}(\text{the white})$

$P(\text{the} | \text{the white}) = \text{freq}(\text{the white the}) / \text{freq}(\text{the white})$

data driven: given $n_1, n_2, n_3, n_4, \dots, n_z$



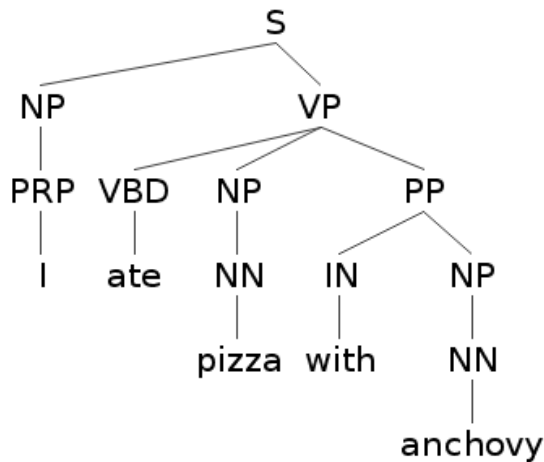
N-gram models

- The higher n , the more context is captured
- The higher n , the less statistical evidence we find for each context: sparse data problem

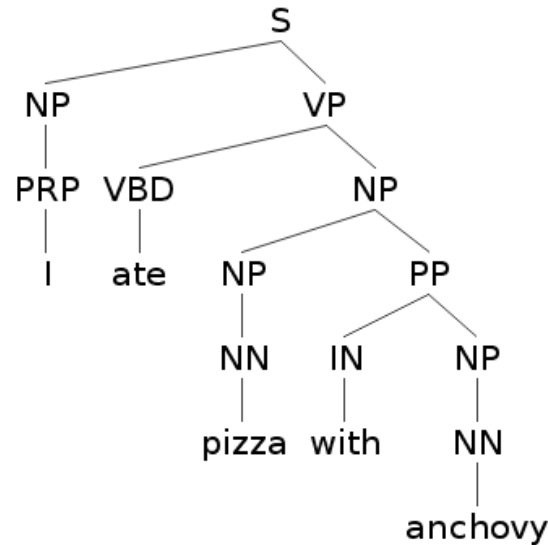
PCFG as a language model



$$P(\text{tree}) = \prod P(\text{rule}_i)$$



16/343



16/2401

$$P(\text{parse}) = \prod p(\text{rule}_i)$$
$$P(\text{sentence}) = \sum p(\text{parse}_k)$$
$$P(\text{text}) = \sum p(\text{sentence}_l)$$



Bayes' Rule & Noisy Channel

	P(Input)	P(Output Input)	
Machine Translation		Translation model	
OCR		Model of OCR errors	
Spellchecking		Language model	Model of spelling errors
POS-Tagging		Tag-Word Model	
Speech Recognition		Acoustic model	



Probabilistic Spelling Correction

Kernighan et al (1990): misspelt word differs from correct word in 1 substitution, insertion, transposition or deletion

error	Correction	Correct Letter	Error Letter	Position	Type
acress	actress	t	-	2	deletion
acress	cress	-	a	0	insertion
acress	caress	ca	ac	0	transposition
acress	access	c	r	2	substitution
...					

- correction = $\operatorname{argmax}_c P(t|c) \cdot P(c)$
with $t = \text{typo}$ and C : list of correct words
- $P(c)$: prior: language model (unigram)
- $P(t|c)$: Model of misspellings



Probabilistic Spelling Correction

- Kernighan: 44×10^6 word AP newswire corpus
- PRIOR:

c	Freq(c)	P(c)
actress	1343	.0000315
cress	0	.000000014
caress	4	.0000001
access	2280	.000058



Probabilistic Spelling Correction

- Kernighan: 44×10^6 word AP newswire corpus
- PRIOR:

c	Freq(c)	P(c)
actress	1343	.0000315
ress	0	.0000000014
caress	4	.00000001
access	2280	.000058

smoothing



- What if we want to calculate the probability of something we haven't seen yet?

I like failblog

+ 'failblog' may not have been seen yet

→ 0 probability

→ 0 probability for entire sentence ($\prod P(\text{word})$)

- Add-1 smoothing: $P(\text{word}) = \frac{\text{freq}(\text{word}) + \alpha}{N + \alpha \cdot d}$

$N + \alpha \cdot d$

with α : normalization factor (often $\alpha=1$)

with N : total number of tokens (words)

with d : total number of types (individual words)



Probabilistic Spelling Correction

- Model of misspellings: $P(t|c)$
- Proper $P(t|c)$ cannot be computed, but can be estimated
- Use corpus of errors to construct confusion matrix of 26×26 for each type of mistake

$\text{del}[x,y]$: count how many times xy was typed as x

$\text{ins}[x,y]$: count how many times x was typed as xy

$\text{sub}[x,y]$: count how many times x was typed as y

$\text{trans}[x,y]$: how many times xy was typed as yx



Probabilistic Spelling Correction

Correction	$P(c)$	$P(t c)$	$p(t c)p(c)$
actress	.0000315	.000117	3.69×10^{-9}
gress	.000000014	.00000144	2.02×10^{-14}
caress	.0000001	.0000164	1.64×10^{-13}
access	.000058	.000000209	1.21×10^{-11}

- access is rewritten as 'actress'
- use more intelligent prior to improve results in context



Bayes' Rule & Noisy Channel

	P(Input)	P(Output Input)
Machine Translation		Translation model
OCR		Model of OCR errors
Spellchecking		Model of spelling errors
POS-Tagging	Language Model	Tag-Word Model
Speech Recognition		Acoustic model



Probabilistic n-gram POS Tagging

- Requires annotated corpus

can/**md** the/**dt** tag/**nn** be/**vb** better/**jjr**

- Unigram: $P(\text{word}|\text{tag})P(\text{tag})$

frequency of the tag for this word in corpus

- Bigram: $P(\text{word}_i|\text{tag}_i) P(\text{tag}_i|\text{tag}_{i-1})$

frequency of the tag for this word in corpus, given previous tag

- Trigram: $P(\text{word}_i|\text{tag}_i) P(\text{tag}_i|\text{tag}_{i-1},\text{tag}_{i-2})$

frequency of the tag for this word in corpus,
given previous two tags

- Good Results, but possible data sparseness problems



Bayes' Rule & Noisy Channel

	P(Input)	P(Output Input)
Machine Translation	Language Model	Translation model
OCR		Model of OCR errors
Spellchecking		Model of spelling errors
POS-Tagging		Tag-Word Model
Speech Recognition		Acoustic model



Modeling English Pronunciation Variation

- Differences in pronunciation
- 2 classes:
 - allophonic variation (due to context)

about	-	[ax b aw]	32%
	-	[ax b aw t]	16%
	-	[ix b aw]	8%

- Lexical variation

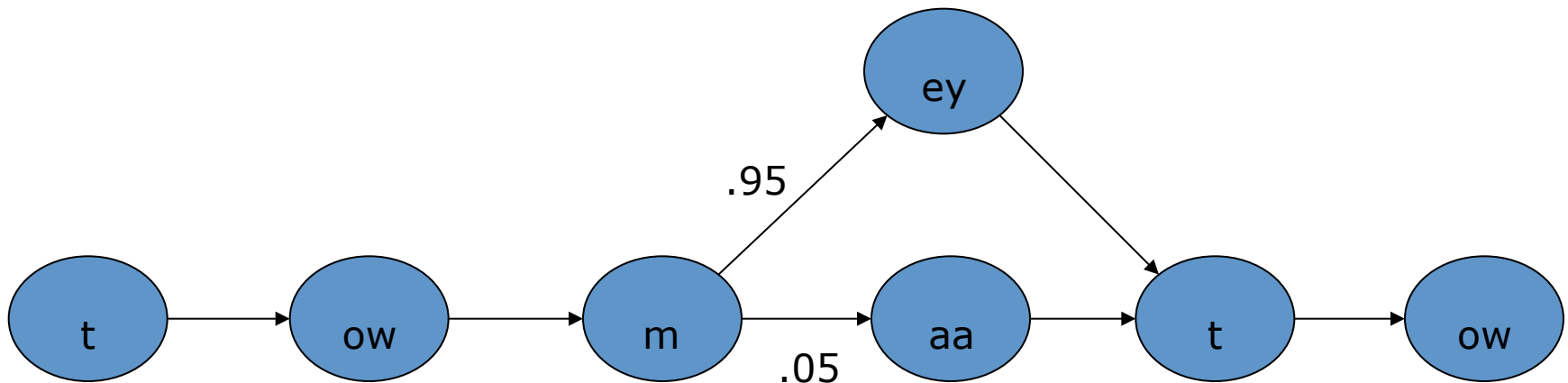
about	-	[baw]	9%
-------	---	-------	----



Modeling English Pronunciation Variation

- we can model the distribution of this variation by introducing probabilities into a FSA

= a Weighted Automaton (Markov Chain)

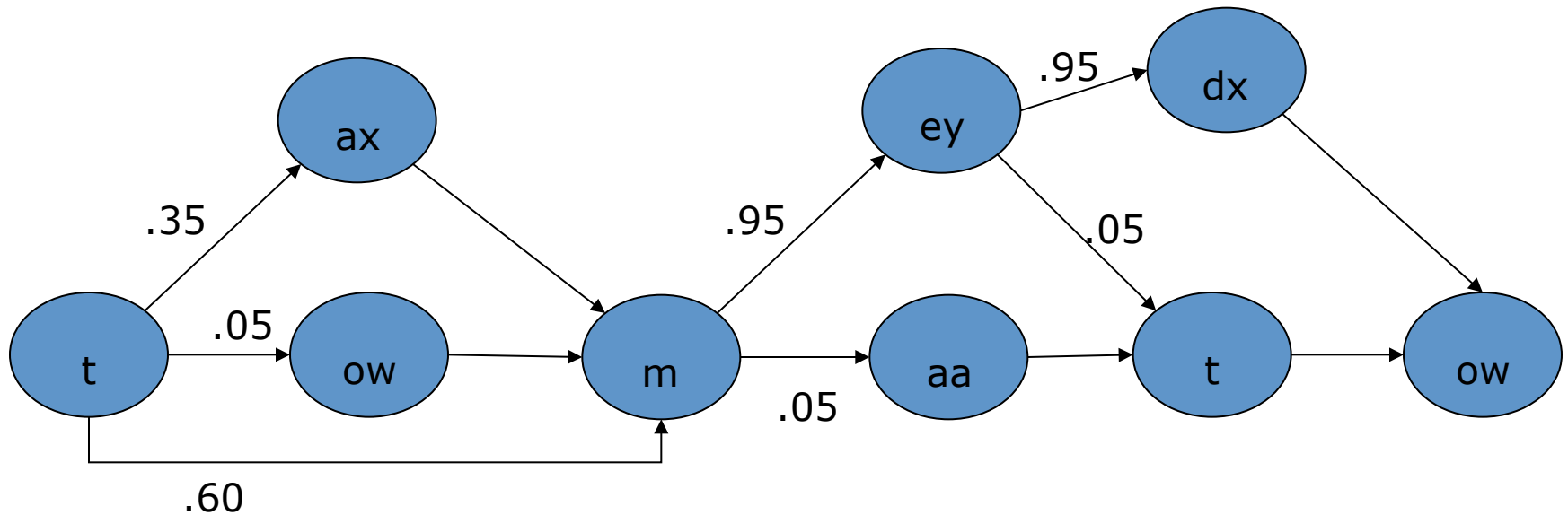


- Models Sociolinguistic variation



Modeling English Pronunciation Variation

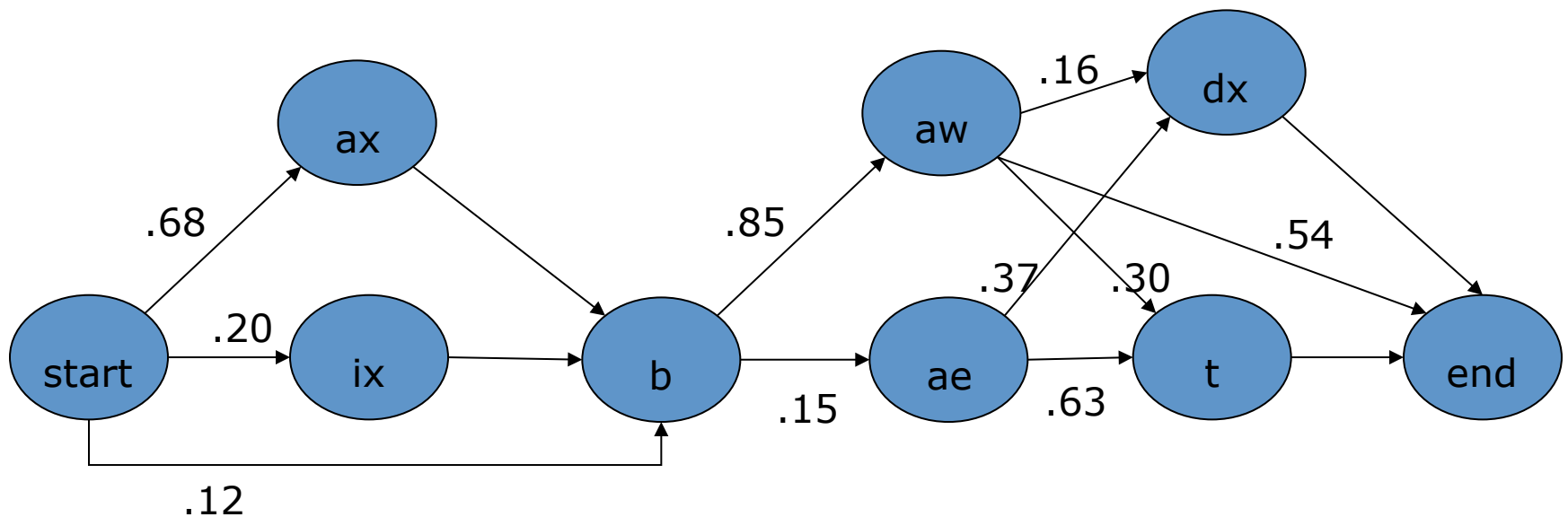
- model allophonic variation:





Modeling English Pronunciation Variation

- “about”: actual weighted automaton trained on pronunciations of Switchboard Corpus





Edit Distance



Minimum Edit Distance

- Spell checking: check writing against list of words/morphotactics
- Suggest list of alternatives?
 - closest match*
 - fuzzy match*
- *How to calculate the "distance" between two words:* minimum edit distance
- The minimal number of deletions, insertions, substitutions to go from word *a* to *b*



Minimum Edit Distance

- *intention* vs *execution*
- 3 operations (deletion, insertion, substitution)
- Alignment

i	n	t	e	*	n	t	i	o	n
*	e	x	e	c	u	t	i	o	n
d	s	s		i	s				

- Levenshtein distance: equal *weight* to all operations, no substitution (1 substitution = 1 deletion + 1 insertion)
- Levenshtein distance of 8 in example above



Minimum Edit Distance

[delete i]

[substitute n for e]

[substitute t for x]

[insert c]

[substitute n for u]

i n t e n t i o n

n t e n t i o n

e t e n t i o n

e x e n t i o n

e x e c n t i o n

e x e c u t i o n



Minimum Edit Distance

- Computed through dynamic programming
- Solve problem by combining solutions to subproblems
- Table-driven

- Useful for
 - Alignment
 - Fuzzy string match
 - Spelling correction



Algorithm

Function LEVENSHTEIN-DISTANCE(*target,source*) returns *levenshtein-distance*

$n \leftarrow \text{length}(\textit{target})$

$m \leftarrow \text{length}(\textit{source})$

Create a distance matrix $\textit{distance}[n+1,m+1]$

Initialize the 0th row and column to be the distance from the empty string

$\textit{distance}[0,0] = 0$

for each column i from 1 to n do

$\textit{distance}[i,0] \leftarrow \textit{distance}[i-1,0] + 1$ (= insertion cost)

for each row j from 1 to m do

$\textit{distance}[0,j] \leftarrow \textit{distance}[0,j-1] + 1$ (= deletion cost)

For each column i from 1 to n do

for each row j from 1 to m do

$\textit{distance}[i,j] \leftarrow \text{MIN}(\textit{distance}[i-1,j] + 1$ (= insertion-cost),
 $\textit{distance}[i,j-1] + 1$ (= deletion-cost),
 $\textit{distance}[i-1,j-1] + 2$ (substitution cost if $A \neq B$)
)

Return $\textit{distance}[n,m]$



Edit distance matrix

n	9	↓8	↙←↓9	↙←↓10	↙←↓11	↙←↓12	↓11	↓10	↓9	↙8
o	8	↓7	↙←↓8	↙←↓9	↙←↓10	↙←↓11	↓10	↓9	↙8	↓9
i	7	↓6	↙←↓7	↙←↓8	↙←↓9	↙←↓10	↓9	↙8	←9	←10
t	6	↓5	↙←↓6	↙←↓7	↙←↓8	↙←↓9	↙8	←9	←10	←↓11
n	5	↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↙←↓9	↙←↓10	↙←↓11	↙↓10
e	4	↙3	←4	↙←↓5	←6	←7	↙←↓8	↙←↓9	↙←↓10	↓9
t	3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↙7	←↓8	↙←↓9	↓8
n	2	↙←↓3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↓7	↙←↓8	↙7
i	1	↙←↓2	↙←↓3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙6	←7	←8
ε	0	1	2	3	4	5	6	7	8	9
	ε	e	x	e	c	u	t	i	o	n

n	9	↓8	↙←↓9	↙←↓10	↙←↓11	↙←↓12	↓11	↓10	↓9	↖8
o	8	↓7	↙←↓8	↙←↓9	↙←↓10	↙←↓11	↓10	↓9	↖8	↓9
i	7	↓6	↙←↓7	↙←↓8	↙←↓9	↙←↓10	↓9	↖8	←9	←10
t	6	↓5	↙←↓6	↙←↓7	↙←↓8	↙←↓9	↖8	←9	←10	←↓11
n	5	↓4	↙←↓5	↙←↓6	↙←↓7	↖←↓8	↙←↓9	↙←↓10	↙←↓11	↙↓10
e	4	↙3	←4	↖←↓5	←6	←7	↙←↓8	↙←↓9	↙←↓10	↓9
t	3	↙←↓4	↖←↓5	↙←↓6	↙←↓7	↙←↓8	↙7	←↓8	↙←↓9	↓8
n	2	↖←↓3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↓7	↙←↓8	↙7
i	1	↙←↓2	↙←↓3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙6	←7	←8
ε	0	1	2	3	4	5	6	7	8	9
	ε	e	x	e	c	u	t	i	o	n

Edit operations are determined by starting from the top right cell, following the arrows to find a path to cell0. Often, several paths are possible.

Path1

		<i>i n t e n t i o n</i>
COL0	<i>[delete i]</i>	<i>n t e n t i o n</i>
COL1	<i>[substitute n for e]</i>	<i>e t e n t i o n</i>
COL2	<i>[substitute t for x]</i>	<i>e x e n t i o n</i>
COL3		<i>e x e n t i o n</i>
COL4	<i>[insert c]</i>	<i>e x e c n t i o n</i>
COL5	<i>[substitute n for u]</i>	<i>e x e c u t i o n</i>
COL6		<i>e x e c u t i o n</i>
...COL9		<i>e x e c u t i o n</i>

n	9	↓8	↙←↓9	↙←↓10	↙←↓11	↙←↓12	↓11	↓10	↓9	↙8
o	8	↓7	↙←↓8	↙←↓9	↙←↓10	↙←↓11	↓10	↓9	↙8	↓9
i	7	↓6	↙←↓7	↙←↓8	↙←↓9	↙←↓10	↓9	↙8	←9	←10
t	6	↓5	↙←↓6	↙←↓7	↙←↓8	↙←↓9	↙8	←9	←10	←↓11
n	5	↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↙←↓9	↙←↓10	↙←↓11	↙↓10
e	4	↙3	←4	↙←↓5	←6	←7	↙←↓8	↙←↓9	↙←↓10	↓9
t	3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↙7	←↓8	↙←↓9	↓8
n	2	↙←↓3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙←↓8	↓7	↙←↓8	↙7
i	1	↙←↓2	↙←↓3	↙←↓4	↙←↓5	↙←↓6	↙←↓7	↙6	←7	←8
ε	0	1	2	3	4	5	6	7	8	9
	ε	e	x	e	c	u	t	i	o	n

Edit operations are determined by starting from the top right cell, following the arrows to find a path to cell0. Often, several paths are possible.

Path1

COL0a [delete i]

COL0b [delete n]

COL0c [delete t]

COL1

COL2 [insert x]

COL3 [substitute n for e]

COL4 [insert c]

COL5 [insert u]

i n t e n t i o n

n t e n t i o n

t e n t i o n

e n t i o n

e n t i o n

e x n t i o n

e x e t i o n

e x e c t i o n

e x e c u t i o n



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5							
e	4							
l	3							
e	2							
d	1							
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6						
e	4	↙←↓5						
l	3	↙←↓4						
e	2	↙←↓3						
d	1	↙←↓2						
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6	↓5					
e	4	↙←↓5	↙↓4					
l	3	↙←↓4	↓3					
e	2	↙←↓3	↙2					
d	1	↙←↓2	↙←↓3					
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6	↓5	↙←↓6				
e	4	↙←↓5	↙↓4	↙←↓5				
l	3	↙←↓4	↓3	↙←↓4				
e	2	↙←↓3	↙2	←↓3				
d	1	↙←↓2	↙←↓3	↙2				
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6	↓5	↙←↓6	↓5			
e	4	↙←↓5	↙↓4	↙←↓5	↙↓4			
l	3	↙←↓4	↓3	↙←↓4	↓3			
e	2	↙←↓3	↙2	←↓3	↙2			
d	1	↙←↓2	↙←↓3	↙2	←3			
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6	↓5	↙←↓6	↓5	↓4		
e	4	↙←↓5	↙↓4	↙←↓5	↙↓4	↙3		
l	3	↙←↓4	↓3	↙←↓4	↓3	↙←↓4		
e	2	↙←↓3	↙2	←↓3	↙2	↙←3		
d	1	↙←↓2	↙←↓3	↙2	←3	←4		
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6	↓5	↙←↓6	↓5	↓4	↙←↓5	
e	4	↙←↓5	↙↓4	↙←↓5	↙↓4	↙3	←↓4	
l	3	↙←↓4	↓3	↙←↓4	↓3	↙←↓4	↙3	
e	2	↙←↓3	↙2	←↓3	↙2	↙←3	←4	
d	1	↙←↓2	↙←↓3	↙2	←3	←4	←5	
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d



Exercise

Calculate the levenshtein distance between
'delen' en 'gedeeld'

n	5	↙←↓6	↓5	↙←↓6	↓5	↓4	↙←↓5	↙←↓6
e	4	↙←↓5	↙↓4	↙←↓5	↙↓4	↙3	←↓4	↙←↓5
l	3	↙←↓4	↓3	↙←↓4	↓3	↙←↓4	↙3	←4
e	2	↙←↓3	↙2	←↓3	↙2	↙←3	←4	←5
d	1	↙←↓2	↙←↓3	↙2	←3	←4	←5	↙←6
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d

n	5	↙←↓6	↓5	↙←↓6	↓5	↓4	↙←↓5	↙←↓6
e	4	↙←↓5	↙↓4	↙←↓5	↙↓4	↙3	←↓4	↙←↓5
l	3	↙←↓4	↓3	↙←↓4	↓3	↙←↓4	↙3	←4
e	2	↙←↓3	↙2	←↓3	↙2	↙←3	←4	←5
d	1	↙←↓2	↙←↓3	↙2	←3	←4	←5	↙←6
ε	0	1	2	3	4	5	6	7
	ε	g	e	d	e	e	l	d

		<i>d e l e n</i>
<i>COL1</i>	<i>[insert g]</i>	<i>g d e l e n</i>
<i>COL2</i>	<i>[insert e]</i>	<i>g e d e l e n</i>
<i>COL3</i>		<i>g e d e l e n</i>
<i>COL4</i>		<i>g e d e l e n</i>
<i>COL5</i>	<i>[insert e]</i>	<i>g e d e e l e n</i>
<i>COL6</i>		<i>g e d e e l e n</i>
<i>COL6b</i>	<i>[delete e]</i>	<i>g e d e e l n</i>
<i>COL7</i>	<i>[substitute n for d]</i>	<i>g e d e e l d</i>

Zie ook: <http://www.let.rug.nl/kleiweg/lev/>



Assignment

```
>>> from nltk.corpus import brown
>>> corpus = brown.sents()
>>> corpus[4]
```

```
[u'The', u'jury', u'said', u'it', u'did', u'find', u'that', u'many', u'of', u'Georgia's', u'registration', u'and', u'election', u'laws', u'', u'are',
u'outmoded', u'or', u'inadequate', u'and', u'often', u'ambiguous', u'', u'.']
```

Write a script that extracts a trigram language model from this corpus. You can do this in 5 steps:

1. Create a dictionary (trigrams = {}) and add all trigrams in the corpus (key) and their associated count (value).
2. Create a dictionary (bigrams = {}) and add all bigrams in the corpus (key) and their associated count (value).
3. For every key in the trigram dictionary, divide the count by the value of the relevant bigram
4. Your trigram dictionary now contains probabilities
5. (save the dictionary using pickle)

Write a script that computes the probability of a sentence, according to your language model

```
>>> probability(corpus[4]) = <some value>
```

DEADLINE: 22 December 2014

Send python code through e-mail to guy.depauw@uantwerpen.be

Don't hesitate to contact your helpline guy.depauw@uantwerpen.be