

Spellcheck

Acknowledgements : Most slides are based on : J&M's Speech and NLP. I
Other sources are acknowledged separately in individual slides.

Two key problems

- Problem 1: Words have a grammar, just as sentences do. How do we map variants of words to their root forms?
- Problem 2: How do we overcome spelling errors?

Corresponding to each of these problems, there is
Science + Engineering

Some questions

- Why is correcting spelling errors important?
- Why are spelling errors made?
- What are the categories of spelling errors?
- What factors can we exploit to recover from them?

Some questions

- Why is correcting spelling errors important?
 - Human errors in typing (search engines)
 - How rampant are spelling errors?
 - OCR
 - Speech-to-text
- Why are spelling errors made?
 - Homophones
 - No neat mapping between structure and pronunciation of words

Britney Spears

Categories of spelling errors

- Non-words
 - *Seperate* for *separate*
- Words
 - *Dessert* for *desert*
 - *Piece* for *peace*

An important question: Is context available?

Categories of spelling errors

- Non-words
 - *Seperate* for *separate*
- Words
 - *Dessert* for *desert*
 - *Piece* for *peace*

An important question: Is context available?

An alternative classification:

- Typographic errors (homologous errors because of keyboard)
- Cognitive errors

How can we recover from spelling errors?

- Edit distance
- Keyboard
- Pronunciation
- Context
- Syntax
- Source-specific issues
 - OCR : cl and d

Two classes of systems

- Spelling error detection
- Spelling error correction

Typing errors

- Single error mis-spellings
 - Insertion
 - Deletion
 - Substitution
 - Transposition

OCR errors

- Correct:

The quick brown fox jumps over the lazy dog.

- Recognized:

'the q~ick brown foxjurnps over tb l azy dog.

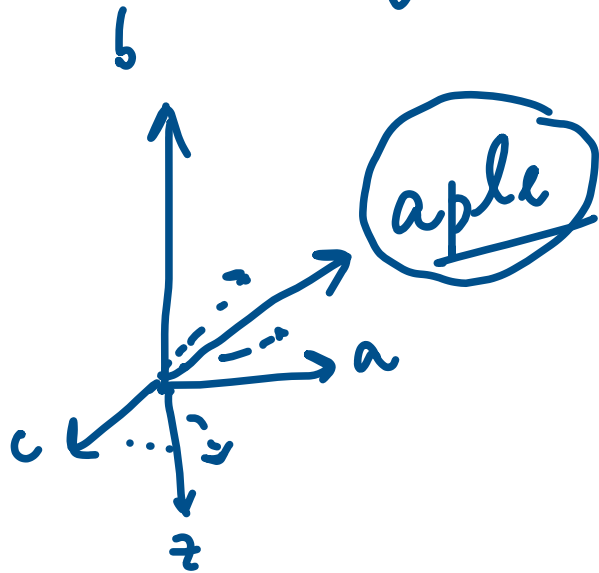
- Substitutions
- Multisubstitutions (framing errors)
- Space deletions
- Space insertions
- Failures

Other issues...

- You may be right, the document may be wrong...
- Google may be unfair to newcomers...

Spelling correction.

aple
✓



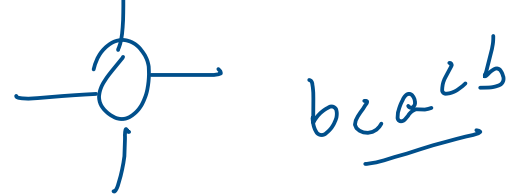
{ apple
 able
 maple
 ape
 ⋮

a b c e l . p z
[1 0 0 1 . . 1 . . 1 , 0 . 0 . .]

→ order _l p

[0 1 0 0 0 4 0 0 30 : 2 . . .]

Order. :



Decision Theoretic ML

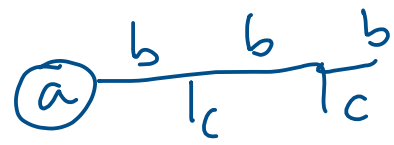
grammar induction.

✓ ML

Syntactic ML

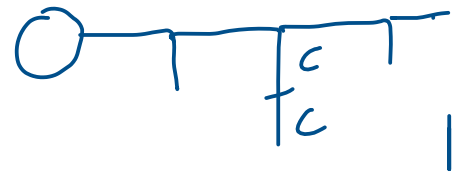
abcbcb

keys



abcbcbcb

○ - a
- b
I - c



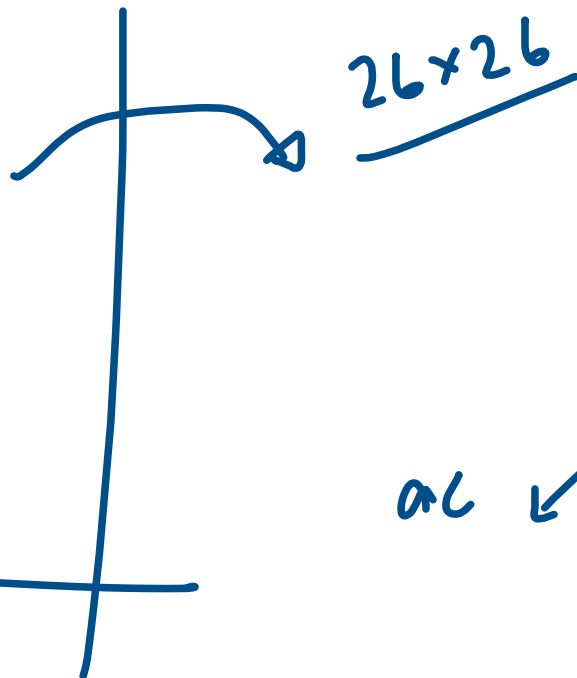
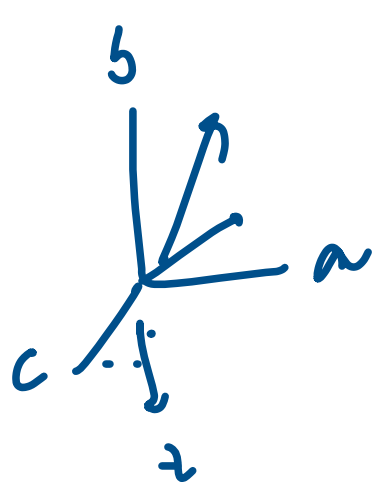
abccbcbcb

keys / no keys

$a(b^+c^+)^+ b^+ c^+$



Syntactic
ML



$$\sqrt{ap_a}$$

$$\{ap, pa\}$$

$$\sqrt{pap}$$

$$\{pa, ap\}$$

$$\frac{aple}{able}$$

$$1, 2,$$

$$\{ap, pl, le\}$$

$$\frac{26 \times 26}{}$$

$$ap \downarrow pl$$

$$aa \dots zz$$

$$[0 \dots i \dots \dots]$$

Bayesian Inference.

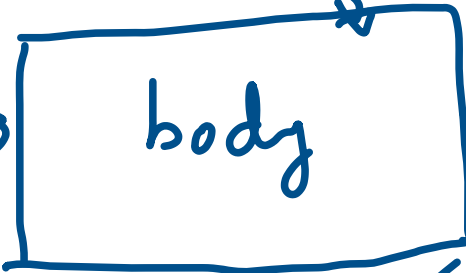
noise

NOISY
CHANNEL

forward

disease

$\begin{bmatrix} D_1 \\ D_2 \\ \vdots \\ D_n \end{bmatrix}$



body

symptom

$\begin{Bmatrix} S_1 \\ \vdots \\ S_m \end{Bmatrix}$

doctor

reverse/inference.

diagnosis. ✓

{ Spellcheck as
Diagnosis }

(MBD) ↑
{ CONSISTENCY BASED
DIAGNOSIS }

(ATE)

Disease

Symptom.

(c)

person
typing

aple type

able
apple
maple

Spellcheck

likelihood

(t)

prior

posterior

$$P(c|t)$$

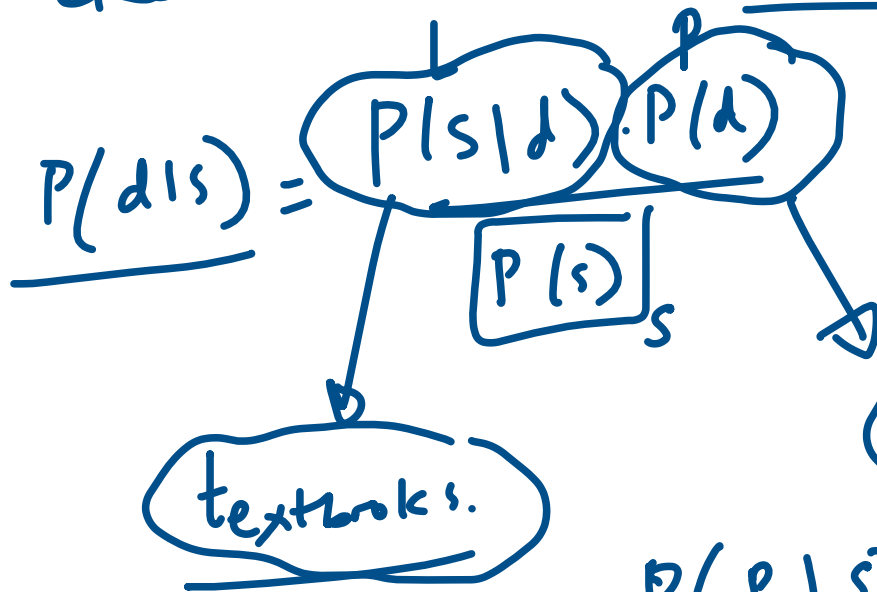
$$= \frac{P(t|c) \cdot P(c)}{P(t)}$$

evidence

$$P(\underset{\substack{\uparrow \\ \text{correct}}}{c} | \underset{\substack{\uparrow \\ \text{type}}}{t}) = \frac{P(\underset{L}{t|c}) \cdot P(\underset{P}{c})}{P(\underset{E}{t})}$$

generative
discriminative

Likelihood \rightarrow Generative process

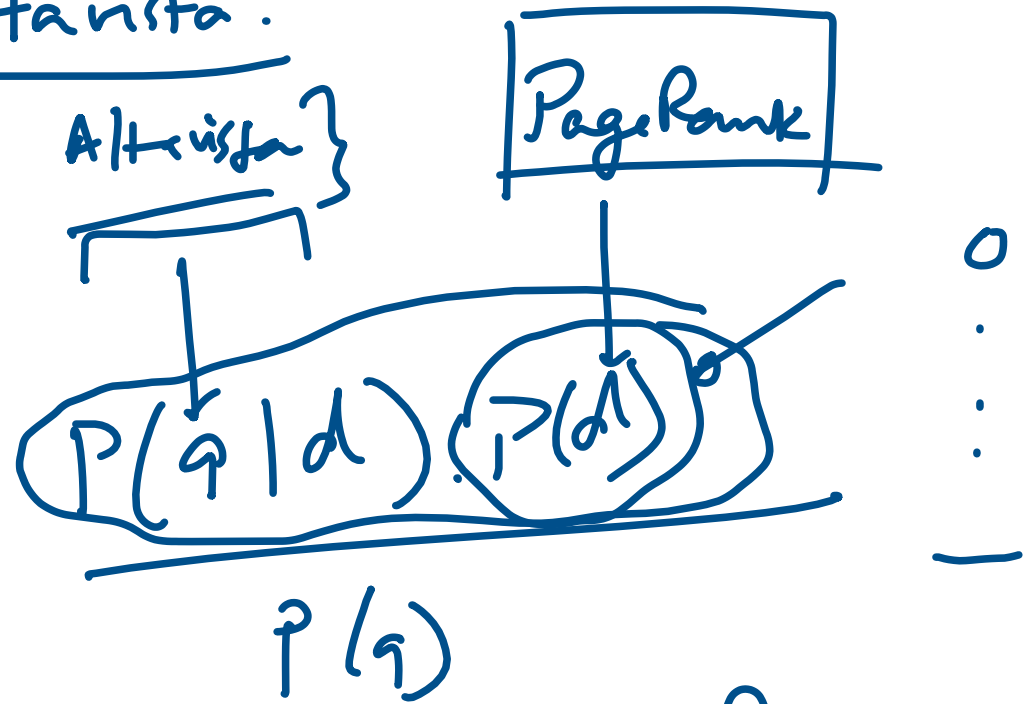


$$\begin{matrix} P(d_1|s) \\ P(d_2|s) \\ \vdots \end{matrix}$$

$$P(\underset{\sim}{p_i} | \underset{\sim}{s}) \propto P(\underset{\sim}{s} | \underset{\sim}{p_i}) \cdot P(\underset{\sim}{p_i})$$

Google vs. Altavista.

PageRank.



$$\frac{P(d|q)}{P(q)} =$$

$$\frac{P(h)}{P(\theta | \underbrace{s}_{H144})} = \frac{P(s|\theta) \cdot P(\theta)}{P(s)}$$

$\frac{P(h)}{P(\theta | \underbrace{s}_{H144})}$

$\frac{P(s|\theta) \cdot P(\theta)}{P(s)}$

$\frac{P(\theta)}{P(s)}$

$\frac{P(s|\theta)}{P(s)}$

$\frac{P(\theta)}{P(s)}$

Subproblems.

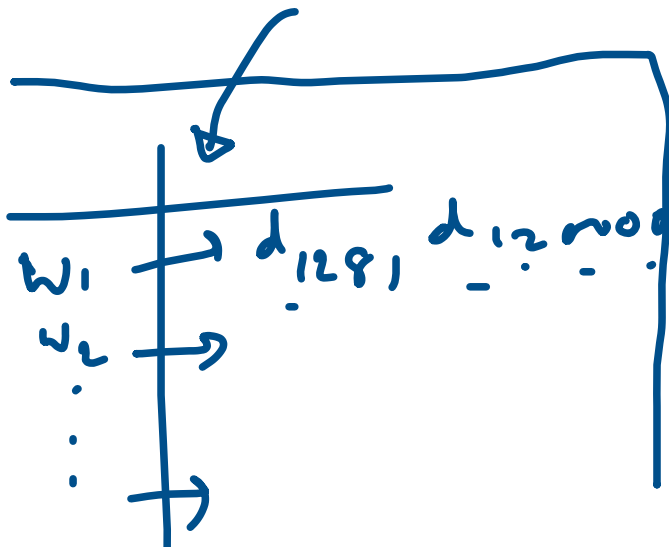
- generate candidates,
- { → prior estimation
- likelihood estimation

Apple → {
able
apple
maple ✓
:
}

bigrams

apple

$\{ap, pl, le\}$



26x26

77

postings file/
inverted index

$aa \rightarrow \sim$
 $ab \rightarrow [abk,]$
.
.
.
.
.
77

Prior

$$P(c|t) \propto P(t|c) \cdot \underbrace{P(c)}_{\text{prior.}}$$

$$\frac{\binom{n}{z}}{n \rightarrow \infty}$$

the
observers

corpus.

{ ...
- a
-
- z

bottom-up.

$$P(c) = \frac{n}{N}$$

$N:$

$a \rightarrow \{-a, \underline{at}, l-\}$

$\{-at-\}$

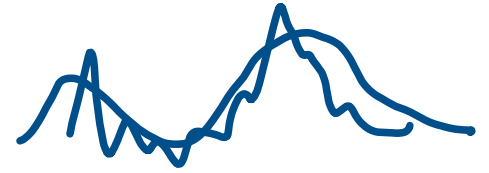
$\{-a, at, t-\}$

[types : no. of distinct words
tokens : no. of words with repetitions.

"the man chase the monkey"

$$\left. \begin{array}{l} V = 4 \\ N = 5 \end{array} \right\}$$

Smoothing.



after smoothing



$$P(c) = \frac{n_c}{N}$$

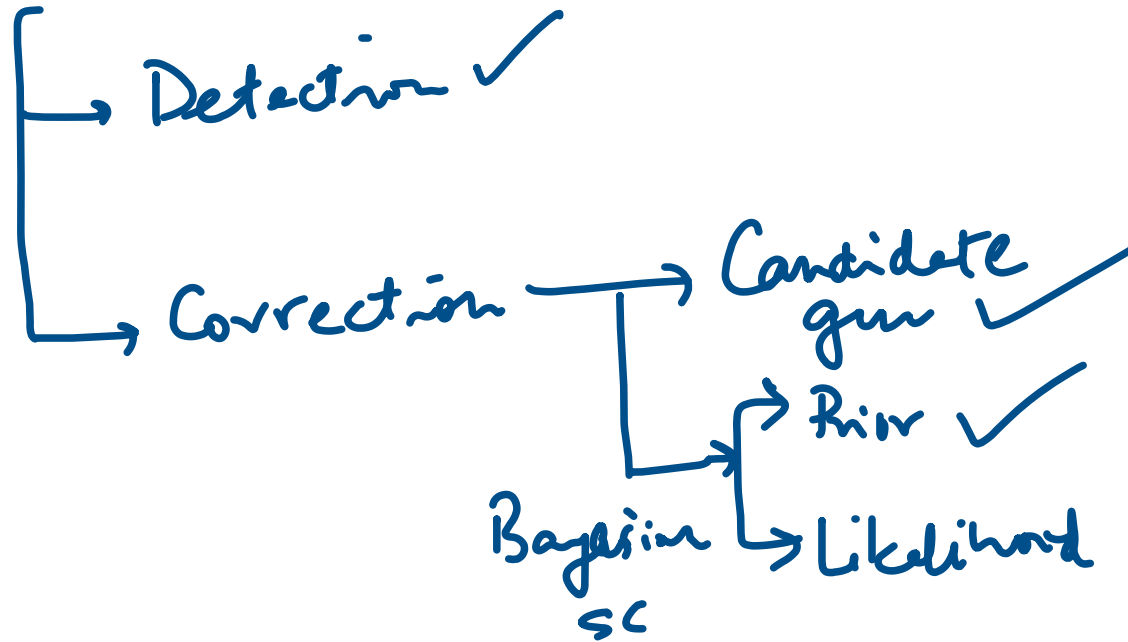
↑

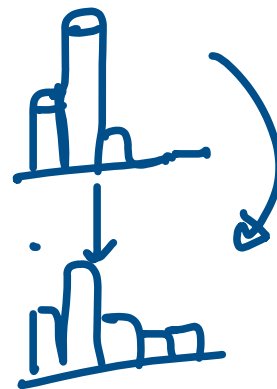
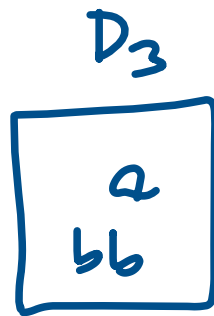
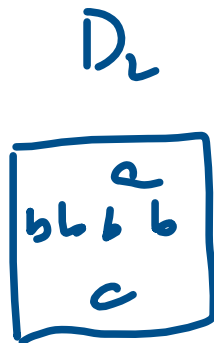
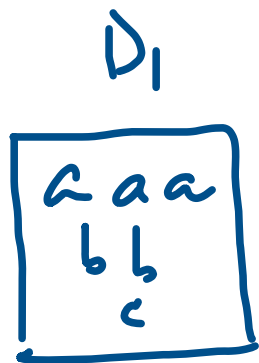
$$\frac{n_c + 1}{N + V}$$

↑.

Revised estimate
of prior
after smoothing

Spellcheck





Laplace (add-1 smoothing)

$$\frac{n_c}{N} \quad \left(\frac{n}{N}\right) \quad \left(\frac{n+1}{N+V}\right) \quad a: \frac{5}{15} \rightarrow \frac{6}{20}$$

$$\frac{n}{N} > \frac{n+1}{N+V}$$

$$nV > N$$

$$n > \frac{N}{V}$$

$$N = 15$$

$$V = 5$$

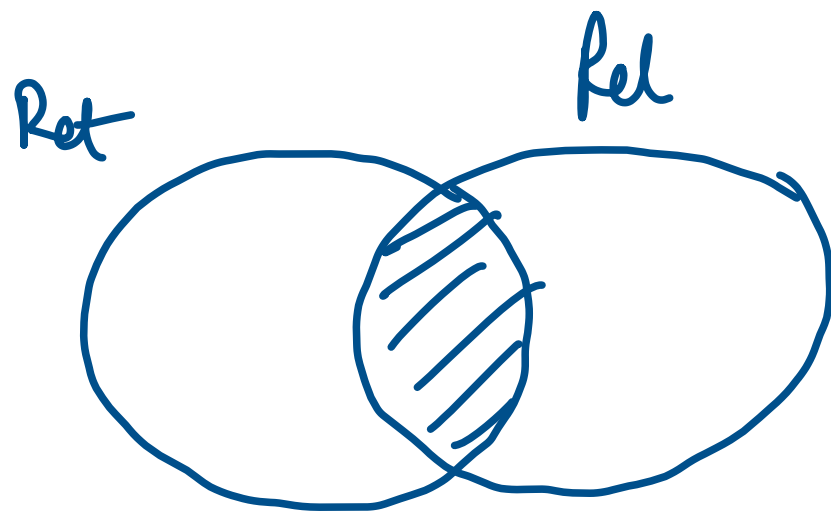
$$n > 3$$

5
types

1.1
tokens

	B.S.	A.S.
a	5	6
b	8	9
c	2	3
d	0	1
e	0	1
	15	20

$$V = 5$$



[Ret: set of retrieved suggestions.]

[Rel: set of relevant suggestions.]

Precision:
$$\frac{|Ret \cap Rel|}{|Ret|} \rightarrow$$

Recall:
$$\frac{|Ret \cap Rel|}{|Rel|} \leftarrow$$

Candidate corrections

t: acress ✓

		Transformation			
Error	Correction	Correct Letter	Error Letter	Position (Letter #)	Type
acress	actress	t	—	2	deletion ✓
acress	cress	—	a	0	insertion
acress	caress	ca	ac	0	transposition
acress	access	c	r	2	substitution
acress	across	o	e	3	substitution
acress	acres	—	2	5	insertion
acress	acres	—	2	4	insertion

c \rightleftarrows t

Prior probabilities

$$\hat{c} = \underset{c \in V}{\operatorname{argmax}} \frac{p(t|c) \cdot p(c)}{L \cdot P}$$

types

c	freq(c)	p(c)
actress	1343	.0000315
cress	0 →	.000000014 →
caress	4	.0000001
access	2280	.000058
across	8436	.00019
acres	2879	.000065

cress →
access

Assoc. press (88) ...

44 million
words

Likelihood.

c: apple
t: apple

t_p

c_{p-1}

→ 'b'

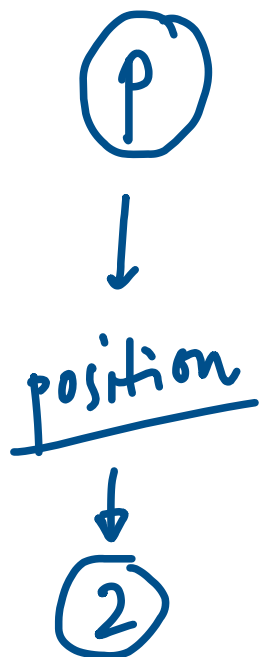
→ 'p'

$$\frac{\text{Count}[pb]}{\text{Count}[p]}$$

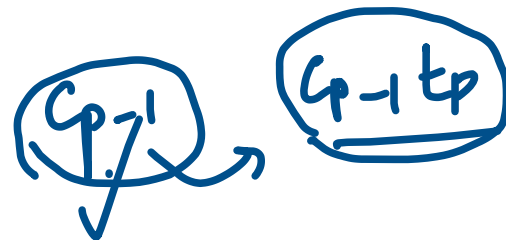
Estimating $p(t|c)$

$\left\{ \begin{array}{l} \text{del}[x,y] : xy \text{ typed as } x \\ \text{ins}[x,y] : x \text{ typed as } \underline{xy} \\ \text{sub}[x,y] : y \text{ was typed as } x \\ \text{trans}[x,y] : xy \text{ typed as } yx \end{array} \right.$

$c \rightarrow t$
Current to ~~type~~



$$P(t|c) = \begin{cases} \frac{\text{del}[c_{p-1}, c_p]}{\text{count}[c_{p-1}c_p]}, & \text{if deletion} \\ \frac{\text{ins}[c_{p-1}, c_p]}{\text{count}[c_{p-1}]}, & \text{if insertion} \\ \frac{\text{sub}[c_p, c_p]}{\text{count}[c_p]}, & \text{if substitution} \\ \frac{\text{trans}[c_p, c_{p+1}]}{\text{count}[c_p c_{p+1}]}, & \text{if transposition} \end{cases}$$



sub[X, Y] = Substitution of X (incorrect) for Y (correct)

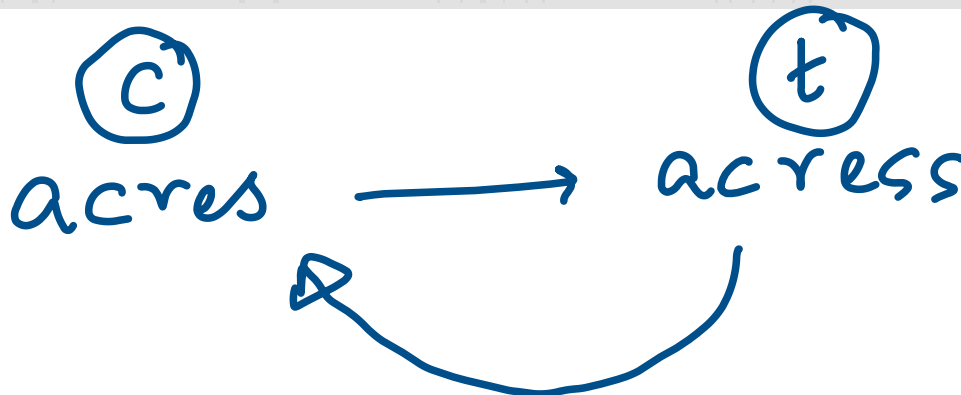
X	Y (correct)																									
	a	b	c	d	e	f	g	h	i	j	k	l	m	n	o	p	q	r	s	t	u	v	w	x	y	z
a	0	0	7	1	342	0	0	2	118	0	1	0	0	3	76	0	0	1	35	9	9	0	1	0	5	0
b	0	0	9	9	2	2	3	1	0	0	0	5	11	5	0	10	0	0	2	1	0	0	8	0	0	0
c	6	5	0	16	0	9	5	0	0	0	1	0	7	9	1	10	2	5	39	40	1	3	7	1	1	0
d	1	10	13	0	12	0	5	5	0	0	2	3	7	3	0	1	0	43	30	22	0	0	4	0	2	0
e	388	0	3	11	0	2	2	0	89	0	0	3	0	5	93	0	0	14	12	6	15	0	1	0	18	0
f	0	15	0	3	1	0	5	2	0	0	0	3	4	1	0	0	0	6	4	12	0	0	2	0	0	0
g	4	1	11	11	9	2	0	0	0	1	1	3	0	0	2	1	3	5	13	21	0	0	1	0	3	0
h	1	8	0	3	0	0	0	0	0	0	2	0	12	14	2	3	0	3	1	11	0	0	2	0	0	0
i	103	0	0	0	146	0	1	0	0	0	0	6	0	0	49	0	0	0	2	1	47	0	2	1	15	0
j	0	1	1	9	0	0	1	0	0	0	0	2	1	0	0	0	0	0	5	0	0	0	0	0	0	0
k	1	2	8	4	1	1	2	5	0	0	0	0	5	0	2	0	0	0	6	0	0	0	4	0	0	3
l	2	10	1	4	0	4	5	6	13	0	1	0	0	14	2	5	0	11	10	2	0	0	0	0	0	0
m	1	3	7	8	0	2	0	6	0	0	4	4	0	180	0	6	0	0	9	15	13	3	2	2	3	0
n	2	7	6	5	3	0	1	19	1	0	4	35	78	0	0	7	0	28	5	7	0	0	1	2	0	2
o	91	1	1	3	116	0	0	0	25	0	2	0	0	0	0	14	0	2	4	14	39	0	0	0	18	0
p	0	11	1	2	0	6	5	0	2	9	0	2	7	6	15	0	0	1	3	6	0	4	1	0	0	0
q	0	0	1	0	0	0	27	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
r	0	14	0	30	12	2	2	8	2	0	5	8	4	20	1	14	0	0	12	22	4	0	0	1	0	0
s	11	8	27	33	35	4	0	1	0	1	0	27	0	6	1	7	0	14	0	15	0	0	5	3	20	1
t	3	4	9	42	7	5	19	5	0	1	0	14	9	5	5	6	0	11	37	0	0	2	19	0	7	6
u	20	0	0	0	44	0	0	0	64	0	0	0	0	2	43	0	0	4	0	0	0	0	2	0	8	0
v	0	0	7	0	0	3	0	0	0	0	0	1	0	0	1	0	0	0	8	3	0	0	0	0	0	0
w	2	2	1	0	1	0	0	2	0	0	1	0	0	0	0	7	0	6	3	3	1	0	0	0	0	0
x	0	0	0	2	0	0	0	0	0	0	0	0	0	0	0	0	0	0	9	0	0	0	0	0	0	0
y	0	0	2	0	15	0	1	7	15	0	0	0	2	0	6	1	0	7	36	8	5	0	0	1	0	0
z	0	0	0	7	0	0	0	0	0	0	0	7	5	0	0	0	0	2	21	3	0	0	0	0	3	0

Ack: Golding paper, 95

Evaluating Candidates

c	freq(c)	p(c)	p(t c)	p(t c)p(c)	%
actress	1343	.0000315	.000117	3.69×10^{-9}	37%
cress	0	.000000014	.00000144	2.02×10^{-14}	0%
caress	4	.0000001	.00000164	1.64×10^{-13}	0%
access	2280	.000058	.000000209	1.21×10^{-11}	0%
across	8436	.00019	.0000093	1.77×10^{-9}	18%
acres	2879	.000065	.0000321 ✓	2.09×10^{-9}	21%
acres	2879	.000065	.0000342 ✓	2.22×10^{-9}	23%

{ 21%
23% } 44%



Knowledge Source for Confusion matrices

- There are lists available on Wikipedia and from Roger Mitton (<http://www.dcs.bbk.ac.uk/~ROGER/corpora.html>) and Peter Norvig (<http://norvig.com/ngrams/>)

But was that right?

... was called a “stellar and versatile **acress** whose combination of sass and glamour has defined her. . .”.

Reference paper:

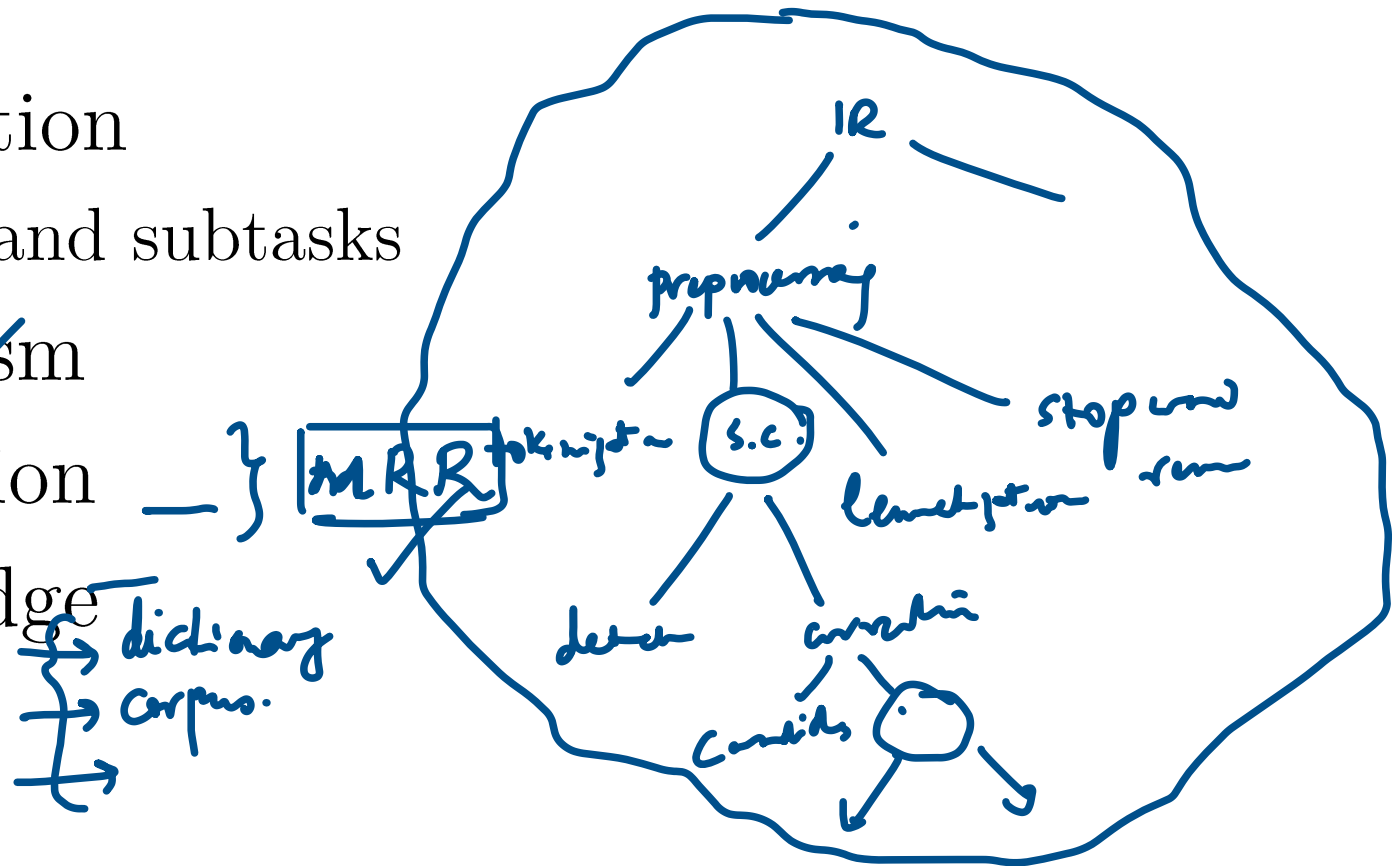
- Kernighan, M. D., Church, K. W. and Gale, W. A. (1990), “***A Spelling Correction Program Based on*** a Noisy Channel Model”, Proceedings of COLING '90, Helsinki

Building a spellcheck application

- Domain
- Application
 - Tasks and subtasks
- Formalism
- Evaluation
- Knowledge

Building an NLP application

- Domain
- Application
 - Tasks and subtasks
- Formalism ✓
- Evaluation
- Knowledge



Stemming/Lemmatization

POS

dancing

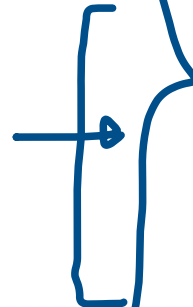


dance

Porter's algorithm.



danc



~ meeting ~