

Natural Language Processing

Matteo Donati
September 29, 2022



Contents

1	Introduction	4
1.1	Levels of Representation	4
2	Basic Text Processing	6
2.1	Regular Expressions	6
2.2	Text Normalization	8
3	Edit Distance, Spelling Correction and the Noisy Channel	9
3.1	Edit Distance	9
3.2	Spelling Correction and the Noisy Channel	10
3.2.1	The Noisy Channel Spelling Method	11
4	N-Gram Language Models	15
4.1	Probabilistic Language Models	15
4.2	Evaluating Models	16
4.3	Smoothing	17
5	Text Classification with Linear Models	19
5.1	Multinomial Naive Bayes Classifier	19
5.1.1	Sentiment Classification	20
5.2	Lexicons for Sentiment, Affect, and Connotation	21
5.3	Logistic Regression	22
5.4	Evaluation Methods	23
6	Vector Semantics and Embeddings - Sparse Representations	24
6.1	Lexical Semantics	24
6.2	Distributional Semantics	24
6.3	Words and Vectors	25
6.3.1	Vector Comparison	26
6.3.2	Reweighting	26
7	Vector Semantics and Embeddings - Dense Representations	28
7.1	Neural Language Modelling	28
7.2	Word2vec	29
7.3	fastText	30
7.4	GloVe	30
7.5	Bias and Embeddings	31
8	Part-of-Speech Tagging	32
8.1	Hidden Markov Model Part-of-Speech Tagging	33
8.2	The Viterbi Algorithm	35

9 Grammars and Parsing	37
9.1 Constituency Grammars	37
9.1.1 Probabilistic Context-Free Grammars	38
9.2 Dependency Grammars	39
10 Sequence Processing with Recurrent Networks	40
10.1 Recurrent Neural Networks	40
10.1.1 Fancier Architectures	43
10.2 Exploding and Vanishing Gradient	44
10.2.1 Long Short-Term Memory	44
10.2.2 Gated Recurrent Units	45
11 Machine Translation, Sequence-to-Sequence and Attention	46
11.1 Machine Translation	46
11.2 Evaluation Metrics	48
11.3 Attention	49
12 Information Extraction and Question Answering	51
12.1 Information Extraction	51
12.2 Question Answering	51
13 CNNs, Transformers and Contextual Word Embeddings	53
13.1 CNNs for NLP	53
13.2 Transformers	53
13.3 Contextual Word Representations	55
14 Dialogue Systems and Chatbots	56
14.1 Dialogue Systems	56
14.2 Chatbots	56
14.3 Task-Oriented Dialogue Agents	57
15 Argument Mining	59
15.1 Data	59
15.2 Methods	59

1 Introduction

Natural Language Processing (NLP) refers to the trend of automating the analysis, generation, and acquisition of human (i.e. natural), language. In particular:

- Analysis (or “understanding”, or “processing”): the input to analysis is language, while the output is some representation that supports useful action.
- Generation: the input to generation is the aforementioned representation, while the output is language.
- Acquisition: is about obtaining the representation and the necessary algorithms from knowledge and data.

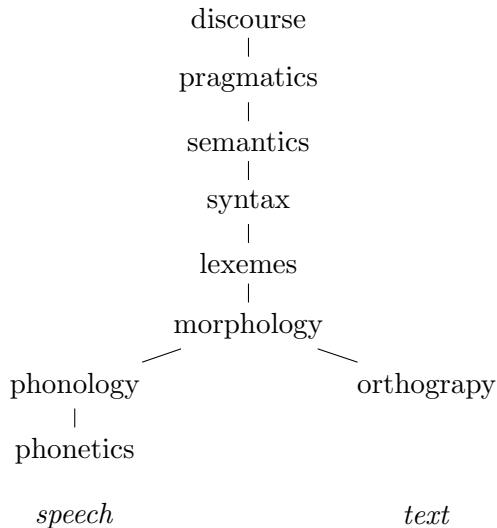
An NLP model is an abstract, theoretical, predictive construct which includes:

- A (partial) representation of the world.
- A method for creating or recognizing worlds.
- A system for reasoning about worlds.

NLP is related with many disciplines, such as computational linguistics, information retrieval, text mining and big data analytics.

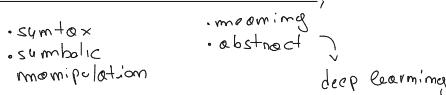
1.1 Levels of Representation

The levels of linguistic representation are the following:



In general, the mappings between levels are extremely complex, and it is difficult to obtain training data for each aspect. Moreover, each string may have many possible interpretations at every level, and the correct resolution of the ambiguity depends on the intended meaning, which is often inferable from the context. In particular:

- **Morphology** is the analysis of words into meaningful components.
- **Lexemen** is about lexical analysis, namely about normalizing and disambiguating words.
- **Syntax** is about transforming a sequence of symbols into a hierarchical or compositional structure. Syntax is closely related to linguistic theories about what makes some sentences well-formed and others not. In general, the number of ambiguities explode combinatorially.
- **Semantics** is about the mapping of natural language sentences into domain representations.
- **Pragmatics** is about any *non-local* meaning phenomena (e.g. “can you pass the salt?” or “Are you 18?”, “Yes, I’m 25”).
- **Discourse** is about texts, dialogues, multi-party conversations.



2 Basic Text Processing

Examples of basic text processing are: regular expressions and text normalization.

2.1 Regular Expressions — formal language for specifying text strings

Regular expressions are used to describe text patterns, to extract information from the text, and to convert text into a convenient/standard form. In particular:

- **Disjunctions** can be of the form:

RE	Match	Example Patterns
/[wW]oodchuck/	Woodchuck or woodchuck	" <u>Woodchuck</u> "
/[abc]/	'a', 'b', or 'c'	"In uomini, in soldati"
/[1234567890]/	any digit	"plenty of <u>7</u> to 5"

Figure 2.2 The use of the brackets [] to specify a disjunction of characters.

disjunction
for single
characters

RE	Match	Example Patterns Matched
/[A-Z]/	an upper case letter	"we should call it ' <u>D</u> rained Blossoms'"
/[a-z]/	a lower case letter	" <u>m</u> y beans were impatient to be hoed!"
/[0-9]/	a single digit	"Chapter <u>1</u> : Down the Rabbit Hole"

Figure 2.3 The use of the brackets [] plus the dash - to specify a range.

RE	Match (single characters)	Example Patterns Matched
/[^A-Z]/	not an upper case letter	"Oyfn pripetchik"
/[^Ss]/	neither 'S' nor 's'	"I <u>h</u> ave no exquisite reason for't"
/[^.]/	not a period	" <u>o</u> ur resident Djinn"
/[e^]/	either 'e' or '^'	"look up <u>^</u> now"
/a^b/	the pattern 'a ^ b'	"look up <u>a ^ b</u> now"

Figure 2.4 The caret ^ for negation or just to mean ^. See below re: the backslash for escaping the period.

disjunction for words

Alternative words can be added using the | pipe (e.g. a|b|c). e.g. [wW]oodchuck s?

- **Optional elements and wildcards** can be of the form: ? for optional previous character, * for zero or more instances of the previous character (Kleene *), + for one or more instances of the previous character (Kleene +), . for any character. e.g. /beg.m/ → anything between "beg" and "m".
 - **Anchors** can be of the form: ^ for the start of a line (e.g. /^The/), \$ for the end of a line (e.g. /The dog\.\$/), b for word boundary (e.g. /\bthe\b/), B for word non-boundary (e.g. /\Bthe\B/).
- dot escape, otherwise ^.
means something else

- More regular expressions operators can be common sets of characters:

RE	Expansion	Match	First Matches
\d	[0-9]	any digit	Party <u>o</u> f <u>5</u>
\D	[^0-9]	any non-digit	Blue <u>m</u> oon
\w	[a-zA-Z0-9_]	any alphanumeric/underscore	Daiyu
\W	[^\w]	a non-alphanumeric	!!!!
\s	[\r\t\n\f]	whitespace (space, tab)	
\S	[^\s]	Non-whitespace	in <u>C</u> oncord

Figure 2.7 Aliases for common sets of characters.

counting operators:

RE	Match
*	zero or more occurrences of the previous char or expression
+	one or more occurrences of the previous char or expression
?	exactly zero or one occurrence of the previous char or expression
{n}	n occurrences of the previous char or expression
{n,m}	from n to m occurrences of the previous char or expression
{n,}	at least n occurrences of the previous char or expression
{,m}	up to m occurrences of the previous char or expression

Figure 2.8 Regular expression operators for counting.

escape symbols:

RE	Match	First Patterns Matched
*	an asterisk “*”	“K*A*P*L*A*N”
\.	a period “.”	“Dr. Livingston, I presume”
\?	a question mark	“Why don’t they come and lend a hand?”
\n	a newline	
\t	a tab	

Figure 2.9 Some characters that need to be backslashed.

Moreover, regular expressions are used in substitutions. In particular, substitutions make use of capture groups to refer to specific sub-parts of the string matching a pattern. For example:

```
/the (.*)er they were, the \1er they will be/
/the (.*)er they (\.*), the \1er we \2/
    |_____
    | capture groups
```

application dependent, these techniques are not always
required

2.2 Text Normalization

Every NLP task needs to do text normalization: segmenting/tokenizing words in running text, normalizing word formats, segmenting sentences in running text. To define some concepts, different examples are considered:

- *I do uh main- / mainly business data processing.* In this sentence there are two kinds of **disfluencies**: fragments and filled pauses.
- *Seuss's cat in the hat is different from other cats!.* *cat* and *cats* are the same **lemma**, but they have different **wordforms**.

Moreover:

- A **type** is an element of a vocabulary, while a **token** is an instance of a type. In general, V is the vocabulary and N the number of tokens. Herdan's law states that $|V| = kN^\beta$. Moreover, **tokenization** could have issues with some specific type of words or part of the text (e.g. contractions, etc.), and with some languages.
- **Word normalization** is the task of putting words/tokens in a standard form (e.g. how to match U.S.A. and USA).
- **Case folding** (e.g. mapping everything to lowercase) is an example of normalization.
- **Lemmatization** reduces inflections or variant forms to their base form (e.g. *am*, *are*, *is* become *be*). An accurate lemmatization requires complete **morphological parsing**.
- **Stemming** is a naive version of morphological analysis which reduces terms to their stems (i.e. the core meaning-bearing units). This implies a crude chopping of affixes (i.e. bits and pieces that adhere to stems). For example, *automate(s)*, *automatic*, *automation* become *automat*.
- **Sentence segmentation** is the process of segmenting sentences in the running text. Question marks and exclamation points are relatively unambiguous, while periods are more ambiguous.

3 Edit Distance, Spelling Correction and the Noisy Channel

3.1 Edit Distance

The concept of edit distance defines how similar are two strings. In particular, the **minimum edit distance** represents the minimum number of edit operations needed to transform one string into another. These edit operations can be: insertions, deletions, substitutions. For example:

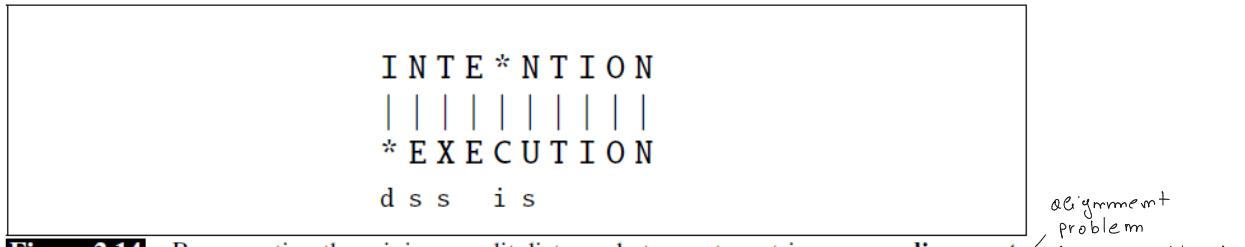


Figure 2.14 Representing the minimum edit distance between two strings as an **alignment**. The final row gives the operation list for converting the top string into the bottom string: d for deletion, s for substitution, i for insertion.

A cost is assigned to each operation. For instance: insertions costs one, deletion costs one, while substitution costs two (but substituting a letter for itself costs zero).

Finding the edit distance can be viewed as a search problem. In particular, one could be interested in finding the shortest path from one string to another. The search space is huge but many paths lead to the same destination. One could exploit dynamic programming so to re-use always visited portions of this space.

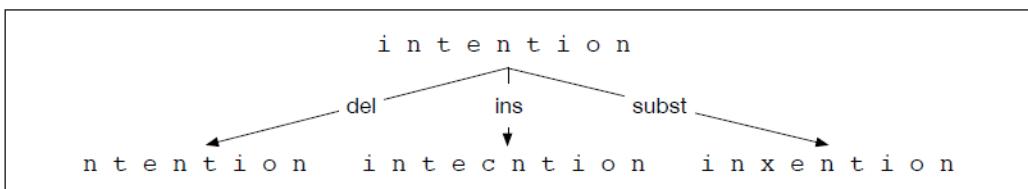


Figure 2.15 Finding the edit distance viewed as a search problem

The algorithm for computing the minimum edit distance is illustrated below:

1. Computation of the minimum edit distance between *intention* and *execution* using Levenshtein distance with cost of one for insertions or deletions, two for substitutions.
2. When entering a value in each cell, one marks which of the three neighbouring cells one came from with up to three arrows. After the table is full one computes an alignment (minimum edit path) by using a backtrace, starting (in this case) at the eight in the lower-right corner and following the arrows back. The sequence of bold cells represents one possible minimum cost alignment between the two strings.

Src \ Tar	#	e	x	e	c	u	t	i	o	n
#	0	1	2	3	4	5	6	7	8	9
i	1	2	3	4	5	6	7	6	7	8
n	2	3	4	5	6	7	8	7	8	7
t	3	4	5	6	7	8	7	8	9	8
e	4	3	4	5	6	7	8	9	10	9
n	5	4	5	6	7	8	9	10	11	10
t	6	5	6	7	8	9	8	9	10	11
i	7	6	7	8	9	10	9	8	9	10
o	8	7	8	9	10	11	10	9	8	9
n	9	8	9	10	11	12	11	10	9	8

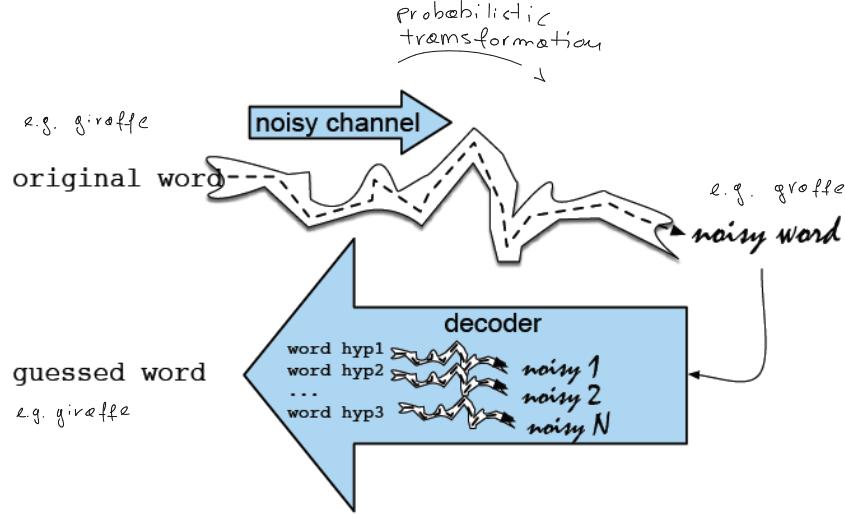
	#	e	x	e	c	u	t	i	o	n
#	0	$\leftarrow 1$	$\leftarrow 2$	$\leftarrow 3$	$\leftarrow 4$	$\leftarrow 5$	$\leftarrow 6$	$\leftarrow 7$	$\leftarrow 8$	$\leftarrow 9$
i	$\uparrow 1$	$\nwarrow 2$	$\nwarrow 3$	$\nwarrow 4$	$\nwarrow 5$	$\nwarrow 6$	$\nwarrow 7$	$\nwarrow 6$	$\leftarrow 7$	$\leftarrow 8$
n	$\uparrow 2$	$\nwarrow 3$	$\nwarrow 4$	$\nwarrow 5$	$\nwarrow 6$	$\nwarrow 7$	$\nwarrow 8$	$\uparrow 7$	$\nwarrow 8$	$\nwarrow 7$
t	$\uparrow 3$	$\nwarrow 4$	$\nwarrow 5$	$\nwarrow 6$	$\nwarrow 7$	$\nwarrow 8$	$\nwarrow 7$	$\leftrightarrow 8$	$\nwarrow 9$	$\uparrow 8$
e	$\uparrow 4$	$\nwarrow 3$	$\leftarrow 4$	$\nwarrow 5$	$\leftarrow 6$	$\leftarrow 7$	$\leftrightarrow 8$	$\nwarrow 9$	$\nwarrow 10$	$\uparrow 9$
n	$\uparrow 5$	$\uparrow 4$	$\nwarrow 5$	$\nwarrow 6$	$\nwarrow 7$	$\nwarrow 8$	$\nwarrow 9$	$\nwarrow 10$	$\nwarrow 11$	$\nwarrow 10$
t	$\uparrow 6$	$\uparrow 5$	$\nwarrow 6$	$\nwarrow 7$	$\nwarrow 8$	$\nwarrow 9$	$\nwarrow 8$	$\leftarrow 9$	$\leftarrow 10$	$\leftarrow 11$
i	$\uparrow 7$	$\uparrow 6$	$\nwarrow 7$	$\nwarrow 8$	$\nwarrow 9$	$\nwarrow 10$	$\uparrow 9$	$\nwarrow 8$	$\leftarrow 9$	$\leftarrow 10$
o	$\uparrow 8$	$\uparrow 7$	$\nwarrow 8$	$\nwarrow 9$	$\nwarrow 10$	$\nwarrow 11$	$\uparrow 10$	$\uparrow 9$	$\nwarrow 8$	$\leftarrow 9$
n	$\uparrow 9$	$\uparrow 8$	$\nwarrow 9$	$\nwarrow 10$	$\nwarrow 11$	$\nwarrow 12$	$\uparrow 11$	$\uparrow 10$	$\uparrow 9$	$\nwarrow 8$

3.2 Spelling Correction and the Noisy Channel

There exist two perspectives to spelling correction:

- **Non-word spelling correction**, for example *graffe* for *giraffe*. Non-word spelling correction can be detected using a dictionary. Then, a ranked list of candidates can be generated, and the ranking is based on a distance metric such as the minimum edit distance. The preference is given to frequent words.
- **Real word spelling correction**, which derives from typographical errors (insertions, deletions, transpositions), or cognitive errors (e.g. *dessert* for *desert*). In real word spelling error detection, any input word could be an error.

In the **noisy channel** model, one can imagine that the surface form one sees is actually a distorted form of an original word passed through a noisy channel. The decoder passes each hypothesis through a model of this channel and picks the word that best matches the surface noisy word:



In this context, the goal is to build a model of the channel. One can pass every word of the language through the channel, so to see which one comes out closest to the observed word. One can see this process as a kind of Bayesian inference: out of all the words in the vocabulary, one tries to find the word w such that $P(w|x)$ is the highest (w is the original word, x is the noisy word):

$$\hat{w} = \underset{w \in V}{\operatorname{argmax}} P(w|x) \quad (1)$$

By Bayes rule:

$$P(w|x) = \frac{P(x|w)P(w)}{P(x)} \quad (2)$$

Because $P(x)$ is independent of w , one has:

$$\hat{w} = \underset{w \in V}{\operatorname{argmax}} P(w|x) = \underset{w \in V}{\operatorname{argmax}} P(x|w)P(w) \quad (3)$$

Finally, one can restrict the analysis to a set of candidates $C \subset V$:

$$\hat{w} = \underset{w \in C}{\operatorname{argmax}} \underbrace{P(x|w)}_{\text{channel model}} \underbrace{P(w)}_{\text{prior}} \quad (4)$$

The prior probability of a word is usually defined by the contest c ($c \rightarrow w \rightarrow x$), so that:

$$\hat{w} = \underset{w \in C}{\operatorname{argmax}} \underbrace{P(x|w)}_{\text{channel model}} \underbrace{P(w|c)}_{\text{language model}} \quad (5)$$

3.2.1 The Noisy Channel Spelling Method

This method can be executed as follows:

1. Find words with similar spelling to the input word. Usually, the majority of error are single-letter changes. To do so, one can use the Damerau-Levenshtein edit distance (includes four edits: insert, delete, substitute, transpose adjacent). Then, a list of candidates at distance one is computed.
2. Score each candidate based on the language model and the channel model. To do so, one could use an $|A| \times |A|$ confusion matrix:

X	sub[X, Y] = Substitution of X (incorrect) for Y (correct)																									
	Y (correct)																									
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	5	0	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	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	2	0	8	0	
v	0	0	7	0	0	3	0	0	0	0	0	1	0	0	1	0	0	8	3	0	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	9	0	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	7	5	0	0	0	0	2	21	3	0	0	0	0	3	0	

Considering the sentence “*was called a stellar and versatile across whose combination of sass and glamour has defined her*”, one can define the following table containing the candidate corrections for the misspelling *acress* and the transformations that would have produced the error:

Error	Correction	Transformation			
		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	—	s	5	insertion
acress	acres	—	s	4	insertion

Supposing priors are computed without taking the context into account, so that $P(w) = \text{count}(w)/|N|$, where N is the number of tokens in the corpus, then (considering COCA, which is composed of 404,253,213 words):

Candidate	Correct	Error	$\frac{\text{channel model}}{\text{(error model)}}$		Prior	$10^9 * P(x w)P(w)$	
			Correction	Letter	x w		
					$P(x w)$	$P(w)$	
actress	t	-	c ct	c	.000117	.0000231	2.7
cress	-	a	a #	a	.00000144	.000000544	0.00078
caress	ca	ac	ac ca	ac	.00000164	.00000170	0.0028
access	c	r	r c	r	.000000209	.0000916	0.019
across	o	e	e o	e	.0000093	.000299	2.8
acres	-	s	es e	es	.0000321	.0000318	1.0
acres	-	s	ss s	ss	.0000342	.0000318	1.0

Figure B.5 Computation of the ranking for each candidate correction, using the language model shown earlier and the error model from Fig. B.4. The final score is multiplied by 10^9 for readability.

One can notice how *actress* is only the second most likely word, and not the first one. Instead, if one also considers the context c so that $P(w_i) = P(w_i|w_{i-1})$:

- $P(\text{actress}|\text{versatile}) = .000021$
- $P(\text{across}|\text{versatile}) = .000021$
- $P(\text{whose}|\text{actress}) = .001$
- $P(\text{whose}|\text{across}) = .000006$

Then:

- $P(\text{versatile actress whose}) = .000021 * .001 = 210 * 10^{-10}$
- $P(\text{versatile across whose}) = .000021 * .000006 = 1 * 10^{-10}$

Lastly:

- $P(\text{versatile actress whose}|\text{versatile across whose}) = 2.7 * 210 * 10^{-19}$
- $P(\text{versatile across whose}|\text{versatile across whose}) = 2.8 * 10^{-19}$

× When considering real word spelling errors, between 25% and 40% of such spelling errors are valid English words. However, one can still use the noisy channel method by considering the following:

1. Given an input sentence $X = \{x_1, x_2, \dots, x_k, \dots, x_n\}$, generate a set of candidate correction sentences $C(X)$.

2. Pick the sentence with the highest language model probability. Here, the channel model accounts for the possibility that a word is not an error, for example:

$$P(x|w) = \begin{cases} \alpha & \text{if } x = w \\ \frac{1-\alpha}{|C(x)|} & \text{if } x \in C(x) \\ 0 & \text{otherwise} \end{cases} \quad (6)$$

Considering the sentence “*This used to belong to thew queen. They are leaving in about fifteen minuets to go to her house.*”, then the noisy channel method would yield:

x	w	x w	P(x w)	P(w w_{i-2}, w_{i-1})	10⁸P(x w)P(w w_{i-2}, w_{i-1})
thew	the	ew e	0.000007	0.48	333
thew	thew		$\alpha=0.95$	9.95×10^{-8}	9.45
thew	thaw	e a	0.001	2.1×10^{-7}	0.0209
thew	threw	h hr	0.000008	8.9×10^{-7}	0.000713
thew	thwe	ew we	0.000003	5.2×10^{-9}	0.00000156

Figure B.6 The noisy channel model on 5 possible candidates for **thew**, with a Stupid Backoff trigram language model computed from the Google N-gram corpus and the error model from [Norvig \(2009\)](#).

4 N-Gram Language Models

4.1 Probabilistic Language Models

Assigning a probability to a sentence allows one to: spell and grammatical error correction, smart typing, machine translation, augmentative and alternative communication systems, speech recognition (e.g. $P(\text{werewolf}) > P(\text{where wolf})$), summarization, question-answering, and similar. The goal is to compute the probability of a sentence or sequence of words:

$$P(W) = P(w_1, w_2, \dots, w_n) \quad (7)$$

One related task is the one of computing the probability of an upcoming word:

$$P(w_5 | \overbrace{w_1, w_2, w_3, w_4}^{\text{context}}) \quad (8)$$

A model that computes $P(W)$ or $P(w_5|w_1, w_2, w_3, w_4)$ is called a **language model**. Moreover, one can use the notation w_i^n to indicate a sequence of n words: w_1, w_2, \dots, w_n . Remembering the chain rule:

$$\begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \cdots P(w_n|w_1^{n-1}) \\ &= \prod_{i=1}^n P(w_i|w_1^{i-1}) \end{aligned} \quad (9)$$

For example: $P(\text{its water is so transparent}) = P(\text{its}) * P(\text{water}|\text{its}) * P(\text{is}|\text{its, water}) * P(\text{so}|\text{its, water, is}) * P(\text{transparent}|\text{its, water, is, so})$. In particular, one could be tempted to estimate such probability as follows:

$$P(\text{the}|\text{its water is so transparent that}) = \frac{\text{count}(\text{its water is so transparent that the})}{\text{count}(\text{its water is so transparent that})} \quad (10)$$

However, there are too many possible sentences and one will never see enough data for estimating these. Thus, one can make the assumption that:

$$P(\text{the}|\text{its water is so transparent that}) \approx P(\text{the}|\underbrace{\text{that}}_{k=1}) \quad (11)$$

or

$$P(\text{the}|\text{its water is so transparent that}) \approx P(\text{the}|\underbrace{\text{transparent, that}}_{k=2}) \quad (12)$$

In general, for $k \geq 1$:

$$P(w_i|w_1^{i-1}) \approx P(w_i|w_{i-k}^{i-1}) \quad (13)$$

Thus:

probability of the word w_i given all the previous words w_1^{i-1} probability of the word w_i given the previous k words w_{i-k}^{i-1}

probability of a sentence

$$P(w_1^n) \approx \prod_{i=1}^n P(w_i | w_{i-k}^{i-1}) \quad (14)$$

The possible types of model are the following:

- **Unigram model:** history does not matter ($k = 0$):

$$P(w_1^n) \approx \prod_i P(w_i) \quad (15)$$

- **Bigram model:** only the previous word matters ($k = 1$):

*the main
problem of
N-gram
models are
long-range
dependencies*

$$P(w_1^n) \approx \prod_i P(w_i | w_{i-1}) \quad (16)$$

- **N-gram models:** the same concept is extended to trigrams, 4-grams, and so on. In general, this is an insufficient model of language (long-range dependencies).

Given these models, one can estimate probabilities using their maximum likelihood estimate. For example, in the case of bigram probabilities:

$$P(w_i | w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})} \quad (17)$$

Finally, one can do everything in log space so to avoid underflow, and also because adding is faster than multiplying ($\log(\prod_i P_i) = \sum_i \log(P_i)$).

4.2 Evaluating Models

There exist two main approaches to evaluate a language model: extrinsic and intrinsic evaluation:

- **Extrinsic/end-to-end/in-vivo/downstream evaluation** is the only way to know if an improvement in a component is really going to help the task at hand. This approach proceeds as follows:

1. Put each model in a task.
2. Run the task, get an accuracy for model A and for model B.
3. Compare the accuracy for A and B.

This approach is the best for comparing model A and model B. However, running NLP systems end-to-end is often very costly.

- **Intrinsic evaluation** measures the quality of component independently of the application. This type of evaluation produces a metric to evaluate the potential improvements in a language model. Such metric is probability-based and it is called **perplexity**. Usually, this approach gives a bad approximation of performances, so in general it is only useful in pilot experiments.

According to perplexity, the best language model is one that best predicts an unseen test set. Perplexity if the inverse probability of the test set, normalized by the number of words:

$$PP(w_1^N) = P(w_1^N)^{-\frac{1}{N}} = \sqrt[N]{\prod_i^{bigrams} \frac{1}{P(w_i|w_{i-1})}} \quad (18)$$

For example, given a sentence consisting of random digits and a model that assigns $P = 1/10$ to each digit, the perplexity of such model can be computed as:

$$PP(w_1^N) = P(w_1^N)^{-\frac{1}{N}} = \left(\underbrace{\frac{1}{10} * \frac{1}{10} * \dots * \frac{1}{10}}_{N \text{ times}} \right)^{-\frac{1}{N}} = \frac{1}{10}^{-1} = 10 \quad \begin{matrix} \max \text{ perplexity} & (\text{equal to} \\ | & \text{the size of possible} \\ \text{digits}) & \end{matrix} \quad (19)$$

In general, minimizing perplexity is equal to maximizing the test set probability according to the given language model.

Some problems of such language models are the following:

- **Overfitting.** Models work well for prediction if the test corpus is similar to the training corpus, at least if these two share the same genre and dialect. N-grams not seen in the training corpus



- **Out-Of-Vocabulary words (OOV).** Open vocabulary systems have to deal with unknown words. To deal with such words one can add a pseudo-word called <UNK>. The solution to this problem is to turn the problem into closed vocabulary:

- either by choosing a vocabulary and modelling all other words as <UNK>;
- or by modelling all least frequent words as <UNK>.

In general, perplexity is artificially reduced by choosing a small vocabulary. Thus, perplexities have to be compared only across language models with the same vocabulary.

4.3 Smoothing

Smoothing is a technique which can be used to avoid the assignment of zero probability to unseen words or contexts. The idea behind this technique is to shave off a bit of probability mass from some more frequent words and give it to unseen events. Below are some of the most important types of smoothing:

- **Laplace smoothing.** This method adds one to all counts and then re-normalizes to obtain a probability. An example on bigram probability estimates is the following:

$$P(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i)}{\text{count}(w_{i-1})} \quad (20)$$

$$P_{\text{Laplace}}^*(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i) + 1}{\text{count}(w_{i-1}) + V} \quad (21)$$

where the adjusted counts are computed as follows:

$$\text{count}^*(w_i) = (\text{count}(w_i) + 1) \frac{N}{N - V} \quad (22)$$

- **Add-k smoothing.** Is a generalization of the Laplace smoothing:

$$P_{\text{Addk}}^*(w_i|w_{i-1}) = \frac{\text{count}(w_{i-1}, w_i) + k}{\text{count}(w_{i-1}) + kV} \quad (23)$$

Usually k is to be optimized.

- **Backoff.** When considering the problem of finding, for example, $P(w_i|w_{i-2}, w_{i-1})$ with no examples of a particular trigram $w_{i-2}w_{i-1}w_i$, one can estimate such probability using $P(w_i|w_{i-1})$ or simply $P(w_i)$. Namely, one can use less context and back off to a lower-order n-gram. In general, to apply backoff smoothing, one needs additional parameters (α, P^*) to discount the higher-order n-grams and save some probability mass:

$$P_{\text{Backoff}}(w_i|w_{i-n+1}^{i-1}) = \begin{cases} P^*(w_i|w_{i-n+1}^{i-1}) & \text{if } \text{count}(w_{i-n+1}^i) > 0 \\ \alpha(w_{i-n+1}^{i-1}) P_{\text{Backoff}}(w_i|w_{i-n+2}^{i-1}) & \text{otherwise} \end{cases} \quad (24)$$

- **Interpolation.** This technique is a generalization of the backoff smoothing:

$$P(w_i|w_{i-2}w_{i-1}) = \lambda_1 P(w_i|w_{i-2}w_{i-1}) + \lambda_2 P(w_i|w_{i-1}) + \lambda_3 P(w_i) \quad (25)$$

with $\sum_i \lambda_i = 1$.

- **Kneser-Ney.** This method uses absolute discounting and looks at the difference between training set and test set to assign probabilities to **<UNK>**.
- **Stupid backoff.** This method was realised to cope with Google n-grams. It uses backoff smoothing without normalization. Thus, it does not compute probabilities:

$$S(w_i|w_{i-k+1}^{i-1}) = \begin{cases} \frac{\text{count}(w_{i-k+1}^i)}{\text{count}(w_{i-k+1}^{i-1})} & \text{if } \text{count}(w_{i-k+1}^i) > 0 \\ \lambda S(w_i|w_{i-k+2}^{i-1}) & \text{otherwise} \end{cases} \quad (26)$$

5 Text Classification with Linear Models

Text categorization is the task of assigning a label or category to an entire text of document. Practical examples of this are: sentiment analysis/opinion mining, SPAM detection, language id, authorship attribution, subject category classification. Indeed, many NLP tasks can be seen as classification tasks.

A **classification task**, formally, is a task in which, given an input x (often d for document), and a fixed set of output classes $Y = y_1, \dots, y_M$ (often c for class), one is supposed to return a predicted class $y \in Y$ for the specific input x . In general, one could use one of the following approaches:

- A rule-based approach.
- A supervised machine learning approach, in which one uses a training set of N labelled documents, $\{d_i, c_i\}$, and a probabilistic classifier which predicts the output class for the a given test sample.
- A generative approach, in which one builds a model of how a class could generate some data. Usually, such systems compare the test data with some probabilistic models to find the best fit for that specific test data. An example of this approach are naive Bayes models.
- A discriminative approach, in which one learns what features are the most useful, without necessarily learning about the classes. An example of this is logistic regression.

5.1 Multinomial Naive Bayes Classifier

The multinomial naive Bayes classifier is a probabilistic classifier which, given d (a document), it returns the class \hat{c} with the maximum posterior probability:

$$\hat{c} = \operatorname{argmax}_{c \in C} P(c|d) \quad (27)$$

Or, equivalently:

$$\hat{c} = \operatorname{argmax}_{c \in C} \underbrace{P(d|c)}_{\text{likelihood}} \underbrace{P(c)}_{\text{prior}} \quad (28)$$

If one represents d as a set of features f_1, \dots, f_n :

$$\hat{c} = \operatorname{argmax}_{c \in C} \underbrace{P(f_1, \dots, f_n|c)}_{\text{likelihood}} \underbrace{P(c)}_{\text{prior}} \quad (29)$$

Assuming mutual feature independence given the class c :

$$c_{NB} = \operatorname{argmax}_{c \in C} P(c) \prod_i P(f_i|c) = \operatorname{argmax}_{c \in C} P(c) \prod_i P(w_i|c) = \operatorname{argmax}_{c \in C} \log P(c) + \sum_i \log P(w_i|c) \quad (30)$$

In particular, $P(c)$ and $P(w_i|c)$ are learned from frequencies in the data. Indeed, the maximum likelihood estimates are:

$$\hat{P}(c) = \frac{N_c}{N_{doc}} \quad (31)$$

$$\hat{P}_{Laplace}(w_i|c) = \frac{\text{count}(w_i, c) + 1}{(\sum_{v \in V} \text{count}(v, c)) + |V|} \quad (32)$$

↑
since +1 is added to
every word in the
vocabulary

5.1.1 Sentiment Classification

Sentiment analysis is about detecting attitudes. A practical example of sentiment analysis is the following:

Category	Documents
-	just plain boring
-	entirely predictable and lacks energy
-	no surprises and very few laughs
+	very powerful
+	the most fun film of the summer
?	predictable with no fun

where:

$$\frac{P(+|\text{predictable with no fun})}{P(-|\text{predictable with no fun})} = \frac{P(+)\text{P}(\text{predictable}|+)\text{P}(\text{no}|+)\text{P}(\text{fun}|+)}{P(-)\text{P}(\text{predictable}|-)\text{P}(\text{no}|-)\text{P}(\text{fun}|-)}$$

Each one of these probabilities can be computed from frequencies. Moreover, there exist some best practices that are commonly used when dealing with sentiment analysis:

- Word occurrence matters more than frequency. Namely, a single occurrence of one word has the same impact of five occurrences of the same word.
- Negation is dealt with by prepending NOT_ to each word between negation and punctuation. For example:

```
didnt like this movie, but I ...
```

becomes:

```
didnt NOT_like NOT_this NOT_movie, but I ...
```

- When the training data is not sufficient, one can use specialized sentiment lexicons. The idea behind lexicons is to focus on specific words instead of using every word as a feature.

list of words containing info about such words

5.2 Lexicons for Sentiment, Affect, and Connotation

There are two influential theories of emotion:

- Basic emotions, where fixed and atomic units are combined into others. An example of this is discrete the following:

Word	anger	anticipation	disgust	fear	joy	sadness	surprise	trust	positive	negative
reward	0	1	0	0	1	0	1	1	1	0
worry	0	1	0	1	0	1	0	0	0	1
tenderness	0	0	0	0	1	0	0	0	1	0
sweetheart	0	1	0	0	1	1	0	1	1	0
suddenly	0	0	0	0	0	0	1	0	0	0
thirst	0	1	0	0	0	1	1	0	0	0
garbage	0	0	1	0	0	0	0	0	0	1

- multi-class classification task

continuous

- Continuum of emotions in 2/3-dimensional space. An example of this is given by space defined by three different features:
 - Valence**, which describes the pleasantness of the stimulus.
 - Arousal**, which describes the intensity of the emotion provoked by the stimulus.
 - Dominance**, which describes the degree of control exerted by the stimulus.

Some lexicons assign to each word a value on all three affective dimensions. An example of this is the following:

	Valence	Arousal	Dominance
vacation	.840	enraged	.962
delightful	.918	party	.840
whistle	.653	organized	.337
consolation	.408	effortless	.120
torture	.115	napping	.046
			powerful
			.991
			authority
			.935
			saxophone
			.482
			discouraged
			.0090
			weak
			.045

Figure 21.4 Samples of the values of selected words on the three emotional dimensions from Mohammad (2018a).

Lexicons, in general, can be created in different ways: by expert annotators, by crowdsourcing, by semi-supervised induction of labels (i.e. propagating from small set of seed labels), by supervised learning.

5.3 Logistic Regression

Logistic regression is a probabilistic classifier the uses supervised machine learning to directly compute $P(c|d)$, instead of building a model of a class c as Naive Bayes does. The main components of a probabilistic machine learning classifier are the following:

- A feature representation of the input $x_i^{(j)}$.
- A classification function computing \hat{y} via $P(y|x)$.
- An objective function involving minimizing an error.
- An optimization algorithm.

Logistic regression computes the output probability by applying the following procedure:

$$z = \left(\sum_{i=1}^n w_i x_i \right) + b = w \cdot x + b \quad (33)$$

then by applying a logistic function to force the output to be a legal probability:

$$P(y = 1) = \sigma(z) = \frac{1}{1 + e^{-z}} = \frac{1}{1 + e^{-(w \cdot x + b)}} \quad (34)$$

$$P(y = 0) = 1 - P(y = 1) = \frac{e^{-(w \cdot x + b)}}{1 + e^{-(w \cdot x + b)}} \quad (35)$$

and lastly by considering 0.5 as decision boundary. An example of input definition is the following:

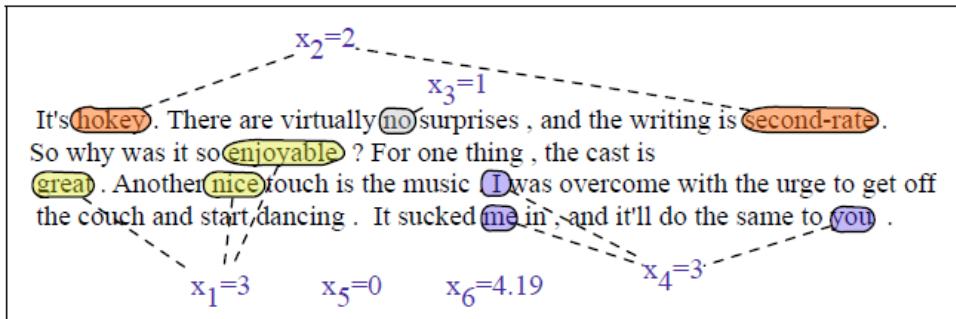


Figure 5.2 A sample mini test document showing the extracted features in the vector x .

where:

- x_1 : count(positive lexicon $\in d$).
- x_2 : count(negative lexicon $\in d$).

- x_3 : 1 if “no” $\in d$, 0 otherwise.
- x_4 : count(1st, 2nd pronouns $\in d$).
— hand-crafted features
- x_5 : 1 if “!” $\in d$, 0 otherwise.
- x_6 : log(word count of d).

The loss function $L(\hat{y}, y)$ is the distance between class and prediction. The main objective is to learn the parameters (w, b) that maximize the probability of the correct labels in the training data given the observations $P(y|x)$:

$$P(y|x) = \hat{y}^y (1 - \hat{y})^{1-y} = \begin{cases} 1 - \hat{y} & \text{if } y = 0 \\ \hat{y} & \text{if } y = 1 \end{cases} \quad (36)$$

$$\log P(y|x) = y \log \hat{y} + (1 - y) \log(1 - \hat{y}) \quad (37)$$

Hence, the cross-entropy loss is obtained by flipping the sign:

$$\text{no local minimum} \quad L_{CE}(\hat{y}, y) = -\log P(y|x) = -[y \log \hat{y} + (1 - y) \log(1 - \hat{y})] \quad (38)$$

L_{CE} for logistic regression is convex. A procedure to find the values for $\theta = (w, b)$ that minimize the loss over training samples (batch of m samples) is given by:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^m \log P(y^{(i)}|x^{(i)}) - \alpha R(\theta) \quad (39)$$

which is exactly what stochastic gradient descent, with mini-batch training, does. The regularization factor $R(\theta)$ is used to avoid overfitting:

- **L2 regularization:** $R(\theta) = \|\theta\|_2^2 = \sum_{j=1}^n \theta_j^2$.
- **L1 regularization:** $R(\theta) = \|\theta\|_1 = \sum_{j=1}^n |\theta_j|$.

In cases which consider more than two classes, the softmax function can be used to normalize the output into probabilities. In this case, the loss function can be computed as:

$$L_{CE}(\hat{y}, y) = - \sum_{c=1}^{|C|} 1\{y = c\} \log P(y = c|x) \quad (40)$$

5.4 Evaluation Methods

In order to evaluate the performance of a machine learning model one can use a gold standard building a contingency matrix and/or a confusion matrix. Having built these matrices, one can either use

- * a micro-average or macro-average type of metrics (where the possible metrics are precision, recall, accuracy, F_1 -score). Moreover, one could also evaluate the model using cross-validation or some statistical significant test.

* micro- F_1 : when data dominated by most frequent class
macro- F_1 : when classes are equally important

6 Vector Semantics and Embeddings - Sparse Representations

6.1 Lexical Semantics

The form of a word is the **lemma** or **citation form** for that word all its derivatives. The **word sense** or **concept** is the meaning component of that particular word. In general, lemmas can be polysemous and the word sense of the word can be used for disambiguation. Moreover:

- **Synonymy** is a relation of (near) identity between two sense of two different words. For example, *couch* and *sofa*. However, the **principle of contrast** states that a difference in linguistic form is always associated with some difference in meaning.
- **Antonymy** is about senses that are opposite with respect with one feature of meaning. For example, *dark* and *light*.
- There exist other relations between senses: **subordination**, **hyponymy** or **IS-A** define the concept of subordinate (e.g. *car* is a subordinate of *vehicle*), **superordination** or **hypernymy** define the opposite of subordination, and **meronymy**.
- Words have few synonyms, but lots of similar words. For example, *cat* and *dog* are not synonyms, but certainly have similar meaning. \ from the meaning view-point
- **Relatedness** is about any relation between words via a semantic frame or field. For example, *car* and *fuel* are somehow related/associated. In general, a **semantic field** is about words that cover a particular domain and bear structured relations with each other. This concept is related to **topic models**. A **semantic frame**, instead, is about words that denote perspectives or participants in a particular type of event.

6.2 Distributional Semantics

Vector semantics define an alternative point of view where words are defined by their usage. In particular, a word is defined by its environment or distribution in the language use, where a distribution is the set of contexts in which it occurs (i.e. the neighbouring words or grammatical environments). Very similar distributions usually imply same meaning. This intuition is called **distributionalist intuition**. However, there exist another intuition, the **vector intuition**, according to which the meaning of a word is represented as a point in an N -dimensional space:



The vector intuition is about embedding each word into a space where similar words are nearby in space. This is the standard way to present meaning in NLP. In particular, to apply such embedding, one would need to choose many things:

- **Matrix design.** Examples are: word \times word, word \times document, word \times discourse context, and others.
- **Reweighting.** Examples are: probabilities, length normalization, TF-IDF, PMI, positive PMI, and others.
- **Dimensionality reduction.** Examples are: latent semantic analysis, non-negative matrix factorization, probabilistic LSA, and others.
- **Vector comparison.** Examples are: Euclidean distance, cosine distance, and others.

choices
to be made
to use
embeddings

6.3 Words and Vectors

A **co-occurrence matrix** represents how often words co-occur. For example:

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

Figure 6.2 The term-document matrix for four words in four Shakespeare plays. Each cell contains the number of times the (row) word occurs in the (column) document.

	As You Like It	Twelfth Night	Julius Caesar	Henry V
battle	1	0	7	13
good	114	80	62	89
fool	36	58	1	4
wit	20	15	2	3

count vector.
vector representation
of a document

Figure 6.3 The term-document matrix for four words in four Shakespeare plays. The red boxes show that each document is represented as a column vector of length four.

is a $|V| \times |D|$ matrix, in which each document is a count vector, identifying a point in a V -dimensional space. Moreover, words can be vectors too. Indeed, one could define a $|V| \times |V|$ **word-word matrix** or **term-term matrix** in which rows are labelled by (target) words, and columns are labelled by (context) words, rather than documents. In the following example, the context is four words on each side:

is traditionally followed by	cherry	pie, a traditional dessert
often mixed, such as	strawberry	rhubarb pie. Apple pie
computer peripherals and personal	digital	assistants. These devices usually
a computer. This includes	information	available on the internet

However, keeping words after the most frequent ones may not be helpful. There is the need to efficiently handle sparse representations.

	context word							
	aardvark	...	computer	data	result	pie	sugar	...
target word	cherry	0	...	2	8	9	442	25
	strawberry	0	...	0	0	1	60	19
	digital	0	...	1670	1683	85	5	4
	information	0	...	3325	3982	378	5	13

Figure 6.5 Co-occurrence vectors for four words in the Wikipedia corpus, showing six of the dimensions (hand-picked for pedagogical purposes). The vector for *digital* is outlined in red. Note that a real vector would have vastly more dimensions and thus be much sparser.

To do so, one could use larger windows and scaling:

	is	then	followed	by	cherry	pie	,	a	dessert	whose
	4	3	2	1	0	1	2	3	4	5
Window: 3	is	<u>then</u>	<u>followed</u>	<u>by</u>	<u>cherry</u>	<u>pie</u>	..	a	dessert	whose
Scaling: flat	0	1	1	1	1	1	1	1	0	
Scaling: $\frac{1}{n}$	0	$\frac{1}{3}$	$\frac{1}{2}$	1	1	1	$\frac{1}{2}$	$\frac{1}{3}$	0	

6.3.1 Vector Comparison

At this point, vectors can be compared using different methods:

- dot-product(w, v) = $w \cdot v = \sum_{i=1}^n w_i v_i$.
- length(v) = $|v| = \sqrt{\sum_{i=1}^n v_i^2}$.
- cosine(w, v) = $\frac{w \cdot v}{|w||v|}$.

6.3.2 Reweighting

The goal of reweighting is to amplify the important, the trustworthy, the unusual, and to de-emphasize the mundane and the quirky. Using raw counts of words is usually a poor proxy for the above rules. One common way of reweighting is **TF-IDF**. This method works by first computing the term frequency and the inverse document frequency:

- **Term frequency (TF)** is the frequency count, usually log-transformed:

$$tf_{t,d} = \log_{10}(\text{count}(t, d) + 1) \quad (41)$$

Words that occur many times in a document have a high tf .

- Inverse document frequency (IDF):

words like "good" or "the" have low idf_t

$$idf_t = \log_{10} \left(\frac{\text{number of documents}}{df_t} \right) \quad (42)$$

Words that occur only in a few documents have a high idf_t . df_t is the number of documents where t occurs.

Lastly:

$$tf-idf_{t,d} = tf_{t,d} * idf_t \quad (43)$$

An alternative to TF-IDF is **positive pointwise mutual information (PPMI)**, which tries to compute at what extent w_1 and w_2 are associated by determining how much more w_1 and w_2 co-occur than by chance:

*equivalent to:
to what extent w_1 and w_2 are associated*

$$PMI(w, c) = \log_2 \frac{P(w, c)}{P(w)P(c)} \quad (44)$$

where $P(w, c)$ is the probability of co-occurrence, and $P(w)P(c)$ is the probability of chance co-occurrence. Lastly:

focus on similarity, not much in distance => take positive part

$$PPMI(w, c) = \max(PMI(w, c), 0) \quad (45)$$

However, PMI is biased towards infrequent events. Indeed, very rare words have very high PMI values. One can solve this problem by either use add- k smoothing or by giving rare words slightly higher probabilities:

$$PPMI_\alpha(w, c) = \max \left(\log_2 \frac{P(w, c)}{P(w)P_\alpha(c)}, 0 \right) \quad (46)$$

where:

$$P_\alpha(c) = \frac{\text{count}(c)^\alpha}{\sum_c \text{count}(c)^\alpha} \quad (47)$$

\ increasing this probability

Typically, $\alpha = 0.75$. This usually works because $P_\alpha(c) > P(c)$ for rare c .

7 Vector Semantics and Embeddings - Dense Representations

In general, TF-IDF and PPMI vectors are long and sparse (i.e. most elements are zero), while dense vectors are usually short and dense. Indeed, these dense vectors usually do not store explicit counts and may be easier to use as features in machine learning (less weights to tune). Examples of dense embeddings are Word2vec, GloVe, fastText. These three embeddings are built on ideas from neural language modelling.

7.1 Neural Language Modelling

Statistical language modelling suffers from the curse of dimensionality, and a lot of information in the training corpus is usually not used. For example, “*the cat is walking in the bedroom*” should help generalize to make “*a dog running in a room*” almost as likely. The idea of neural language modelling is to predict rather than to count. Each word $w \in V$ is associated with a **word feature vector** $C(w) \in \mathbb{R}^m$, and such vectors are used to express the joint probability function of sequences $f(w_t, \dots, w_{t-n+1}) = \hat{P}(w_t | w_1^{t-1})$, where C and f are learned simultaneously.

$$\hat{P}(w_t = i | w_1^{t-1}) = f(i, w_t, \dots, w_{t-n+1}) = g(i, C(w_{t-1}, \dots, C(w_{t-n+1}))) \quad (48)$$

In general, similar words are expected to have similar feature vectors.

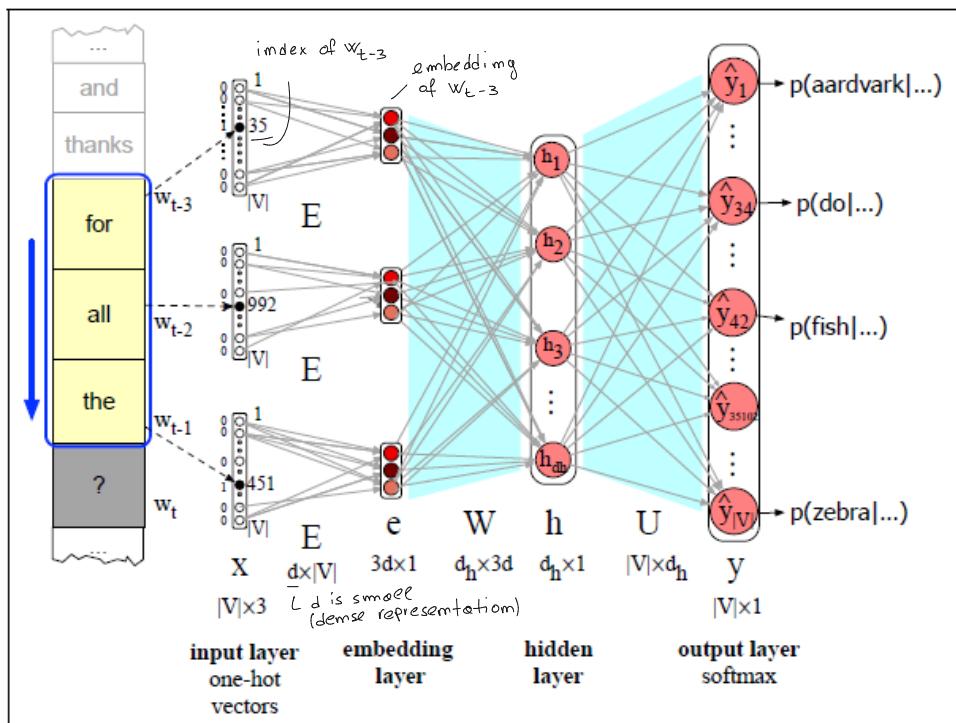


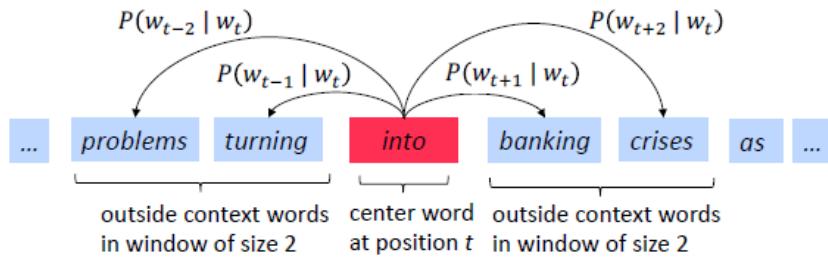
Figure 7.13 Forward inference in a feedforward neural language model.

The main goal is
to learn E , which
can be used to encode
words of, e.g., other
texts

7.2 Word2vec

Word2vec is the most popular embedding methods. It uses the running text as implicitly supervised training data. In particular, instead of counting how often each word w occurs near *apricot*, a classifier is trained to predict whether w is likely to show up near *apricot*. There is no need for hand-labelled supervision, a word s near *apricot* acts as gold standard to the problem. The result of this process will be a set of trained weights which will be used as word embeddings. The idea here is the following:

- Given a large corpus of text, every word in a fixed vocabulary is represented by a vector. One simply needs to iterate through each word in the corpus, fixing the target, or **center word** (c), and the context (o).
- The similarity between c and o is used to calculate the probability of o given c (or vice versa).
- The word vectors are adjusted so to maximize such probability.



Assuming all context words to be independent, for each position $t = 1, \dots, T$, the context words within a window of fixed size m , given the center word w_t , are predicted:

$$Likelihood = L(\theta) = \prod_{t=1}^T \prod_{-m \leq j \leq m, j \neq 0} P(w_{t+j} | w_t; \theta) \quad (49)$$

The objective function, $J(\theta)$, is the average negative log-likelihood:

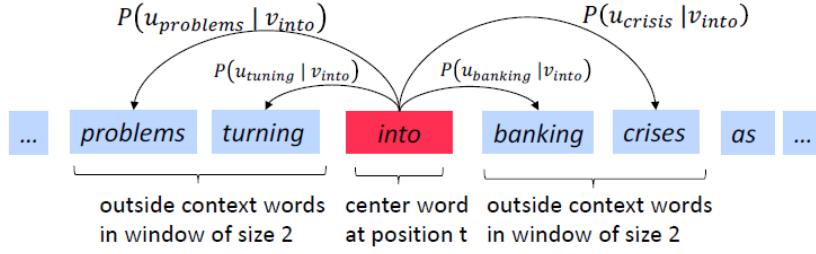
$$J(\theta) = -\frac{1}{T} \log L(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m, j \neq 0} \log P(w_{t+j} | w_t; \theta) \quad (50)$$

To calculate $P(w_{t+j} | w_t; \theta)$ one usually uses two vectors for the same word w : u_w when w is an outside word (u_o), and v_w when w is a center word (v_c). At this point one can compute:

$$P(o|c) = \frac{\exp(u_o^\top v_c)}{\sum_{w \in V} \exp(u_w^\top v_c)} \quad (51)$$

An example of this is the following:

θ are the parameters that enable the embedding. These will store two different vectors for the same word w : one for the case in which w is a target word, and one for the case in which w is a context word.



where $P(u_{problems}|v_{into})$ is a shortcut for $P(\text{problems}|into; u_{problems}, v_{into}, \theta)$.

7.3 fastText

fastText is an extension of Word2vec which uses subword models to deal with OOV words and sparsity in languages with rich morphology. Each word is represented as itself plus a bag of constituent n-grams. For instance, with $n = 3$, *where* is represented by `<where>`, `<wh, whe, her, ere, re>`. Lastly, an embedding is learned for each constituent n-gram, and *where* is represented by the sum of all such embeddings.

7.4 GloVe

In general, there are two families of methods for learning word vectors: global matrix factorization methods, and local context window methods (based on learning the embeddings thanks to neural networks). GloVe combines the ideas of count matrices and neural network predictions. The aim is to have components of meaning as linear operations in the vector space. In particular, it uses **ratios of co-occurrence probabilities** to encode meaning components:

Table 1: Co-occurrence probabilities for target words *ice* and *steam* with selected context words from a 6 billion token corpus. Only in the ratio does noise from non-discriminative words like *water* and *fashion* cancel out, so that large values (much greater than 1) correlate well with properties specific to ice, and small values (much less than 1) correlate well with properties specific to steam.

Probability and Ratio	$k = \text{solid}$	$k = \text{gas}$	$k = \text{water}$	$k = \text{fashion}$
$P(k \text{ice})$	1.9×10^{-4}	6.6×10^{-5}	3.0×10^{-3}	1.7×10^{-5}
$P(k \text{steam})$	2.2×10^{-5}	7.8×10^{-4}	2.2×10^{-3}	1.8×10^{-5}
$P(k \text{ice})/P(k \text{steam})$	8.9	8.5×10^{-2}	1.36	0.96

Annotations on the right side of the table:

- A brace groups the last two columns ($k = \text{water}$ and $k = \text{fashion}$) with the label "similar".
- An arrow points from the "similar" group down to the ratio row, with the label "measuring 'water' and 'fashion' more discriminative so to distinguish the meaning of 'ice' and 'steam'".

These ratios of co-occurrence probabilities can be captured as linear meaning components in a word vector space thanks to a **log-bilinear model**: if the dot product can be made equal to the log of the co-occurrence probability

$$w_i \cdot w_j = \log P(i|j) \quad (52)$$

then a **vector difference** turns into a ratio of the co-occurrence probabilities:

$$w_x \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)} \quad (53)$$

To obtain that, the loss function is defined as:

$$J = \sum_{i,j=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2 \quad (54)$$

7.5 Bias and Embeddings

Examples of embeddings encode cultural stereotypes are the following:

- Embeddings pinpoint sexism implicit in text.
- Embeddings reflect and replicate all sorts of pernicious biases. An example of this is that embeddings for competence adjectives are biased toward men. Fortunately this bias is slowly decreasing.
- Embeddings reflect ethnic stereotypes over time.

8 Part-of-Speech Tagging

Parts of speech/word classes/syntactic categories are useful abstractions which reveal a lot about words and their neighbours. Parts of speech are defined based on syntactic and morphological function. There exist two broad super-categories:

- **Closed class** types, such as prepositions like *from* and *to*, and function words that tend to be short, frequent, and have structuring uses in grammar. Examples of these types are propositions, particles, determiners, conjunctions, pronouns, auxiliary verbs, numerals, and others.
- **Open class** types, such as nouns, verbs, adjectives, and adverbs.

The Penn Treebank POS tags are defined as follows:

Tag	Description	Example	Tag	Description	Example	Tag	Description	Example
CC	coordinating conjunction	<i>and, but, or</i>	PDT	predeterminer	<i>all, both</i>	VBP	verb non-3sg present	<i>eat</i>
CD	cardinal number	<i>one, two</i>	POS	possessive ending	<i>'s</i>	VBZ	verb 3sg pres	<i>eats</i>
DT	determiner	<i>a, the</i>	PRP	personal pronoun	<i>I, you, he</i>	WDT	wh-determ.	<i>which, that</i>
EX	existential 'there'	<i>there</i>	PRP\$	possess. pronoun	<i>your, one's</i>	WP	wh-pronoun	<i>what, who</i>
FW	foreign word	<i>mea culpa</i>	RB	adverb	<i>quickly</i>	WP\$	wh-possess.	<i>whose</i>
IN	preposition/ subordin-conj	<i>of, in, by</i>	RBR	comparative	<i>faster</i>	WRB	wh-adverb	<i>how, where</i>
JJ	adjective	<i>yellow</i>	RBS	superlatv. adverb	<i>fastest</i>	\$	dollar sign	\$
JJR	comparative adj	<i>bigger</i>	RP	particle	<i>up, off</i>	#	pound sign	#
JJS	superlative adj	<i>wildest</i>	SYM	symbol	<i>+, %, &</i>	"	left quote	' or "
LS	list item marker	<i>1, 2, One</i>	TO	"to"	<i>to</i>	"	right quote	' or "
MD	modal	<i>can, should</i>	UH	interjection	<i>ah, oops</i>	(left paren	[, (, {, <
NN	sing or mass noun	<i>llama</i>	VB	verb base form	<i>eat</i>)	right paren],), }, >
NNS	noun, plural	<i>llamas</i>	VBD	verb past tense	<i>ate</i>	,	comma	,
NNP	proper noun, sing.	<i>IBM</i>	VBG	verb gerund	<i>eating</i>	.	sent-end punc	. ! ?
NNPS	proper noun, plu.	<i>Carolinas</i>	VBN	verb past part.	<i>eaten</i>	:	sent-mid punc	: ; ... --

Figure 8.1 Penn Treebank part-of-speech tags (including punctuation).

Part-of-speech tagging is about assigning a POS marker to each input token. This can be seen as a disambiguation task, and the goal of POS tagging is to resolve ambiguities. For example, *that, back, down, put* and *set* are among the most frequent ambiguous words:

- *earning growth took a back/JJ seat.*
- *a small building in the back/NN.*
- *a clear majority of senators back/VBP the bill.*
- *Dave began to back/VB toward the door.*
- *enable the country to buy back/RP its debt.*
- *I was twenty-one back/RB then.*

However, not all tags are equally likely.

× 8.1 Hidden Markov Model Part-of-Speech Tagging

This process relies on a probabilistic sequence model (sequence classifier), to assign a label to each unit in a sequence. Given a sequence of units, this model returns a probability distribution over the possible sequences of labels. Such models are based on the **Markov assumption**, in which only the current state matters:

$$P(q_i = a | q_1 \dots q_{i-1}) = P(q_i = a | q_{i-1}) \quad (55)$$

and are based on the **Markov chain**:

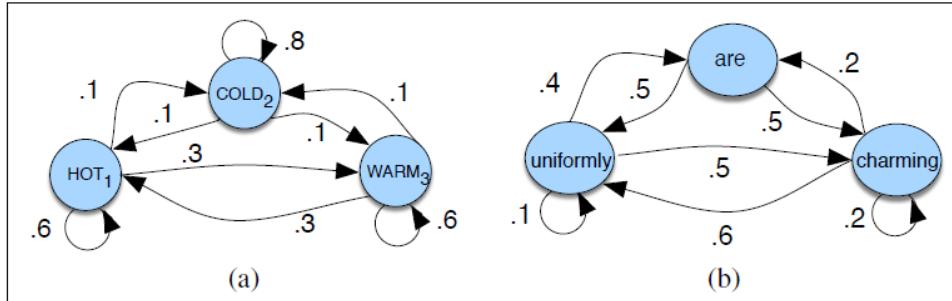


Figure 8.3 A Markov chain for weather (a) and one for words (b), showing states and transitions. A start distribution π is required; setting $\pi = [0.1, 0.7, 0.2]$ for (a) would mean a probability 0.7 of starting in state 2 (cold), probability 0.1 of starting in state 1 (hot), etc.

Formally, a Markov chain is defined by a tuple composed of three elements:

- $Q, \{q_i\}_{i=1}^N$, set of N states.
- $A, [a_{ij}]_{i,j=1}^N$, a transition probability matrix A , where each a_{ij} represents the probability of moving from state i to state j , such that $\sum_{j=1}^N a_{ij} = 1$.
- $\pi, [\pi_i]_{i=1}^N$, an initial probability distribution over states, such that $\sum_{i=1}^N \pi_i = 1$. π_i is the probability that the Markov chain will start in state i .
- $O, [o]_{i=1}^T$, a sequence of T observations, each one drawn from a vocabulary $V = v_1, \dots, v_V$.
- $B, [b_i(o_t)]_{i=1, t=1}^{N, T}$, a sequence of observation likelihoods, also called emission probabilities, each expressing the probability of an observation o_t being generated from a state q_i .

Usually, the events one is interested in are hidden. For instance, POS tags are not explicit in the text. An HMM represents both observed and hidden events. Based on the Markov assumption, then the output independence assumption states that:

$$P(o_i | q_1 \dots q_i \dots q_T, o_1 \dots o_i \dots o_T) = P(o_i | q_i) \quad (56)$$

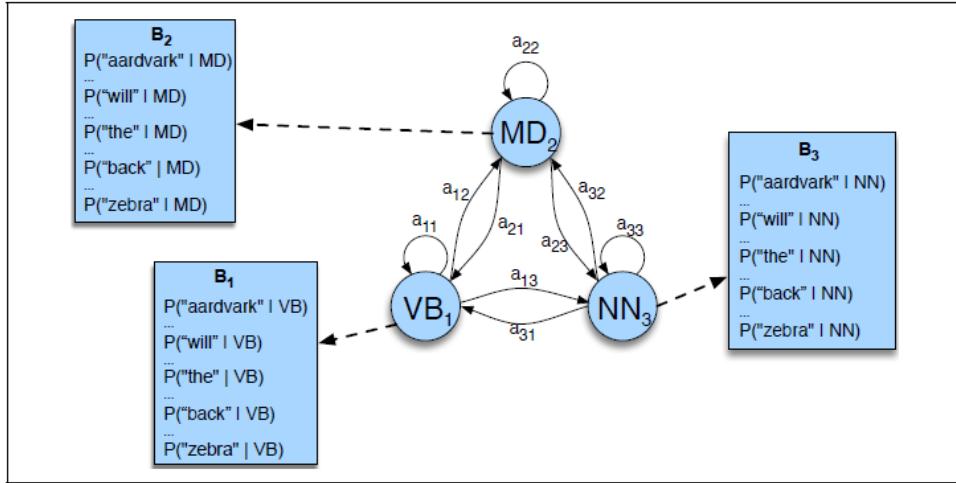


Figure 8.4 An illustration of the two parts of an HMM representation: the A transition probabilities used to compute the prior probability, and the B observation likelihoods that are associated with each state, one likelihood for each possible observation word.

An HMM tagger is composed of two main components: A and B probabilities. For example, *Janet will back the bill*:

- $A : P(t_i|t_{i-1})$:

$$P(VB|MD) = \frac{\text{count}(MD, VB)}{\text{count}(MD)}$$

- $B : P(w_i|t_i)$:

$$P(back|VB) = \frac{\text{count}(VB, back)}{\text{count}(VB)}$$

Decoding is about finding the most probable sequence of hidden states, starting from an observations sequence. Given an HMM $\lambda = (A, B)$, and a sequence of observations $w_1^n = w_1 \dots w_n$, the most probable sequence of hidden states \hat{t}_1^n is given by:

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(t_1^n | w_1^n) \quad (57)$$

By Bayes' rule:

$$\hat{t}_1^n = \underset{t_1^n}{\operatorname{argmax}} P(w_1^n | t_1^n) P(t_1^n) \quad (58)$$

By the Markov assumption (bigram assumption):

$$\hat{t}_1^n \approx \underset{t_1^n}{\operatorname{argmax}} P(w_1^n | t_1^n) \prod_{i=1}^n P(t_i | t_{i-1}) \quad (59)$$

By output independence:

$$\hat{t}_1^n \approx \operatorname{argmax}_{t_1^n} \prod_{i=1}^n \underbrace{P(w_i|t_i)}_{\text{emission}} \underbrace{P(t_i|t_{i-1})}_{\text{transition}} \quad (60)$$

× 8.2 The Viterbi Algorithm

The Viterbi algorithm is a dynamic programming decoding algorithm for HMMs. It uses a probability matrix (lattice) $[v_t(j)]_{1 \leq t \leq T, 1 \leq j \leq N}$, where the columns are observations (tokens), the rows are hidden states (tags), and each cell represents the probability that the HMM is in state j after the first t observations and passing through the most probable state sequence $q_1 \dots q_{t-1}$. Moreover, $v_t(j)$ is recursively defined based on the most probable path leading to it:

$$v_t(j) = \max_{q_1^{t-1}} P(q_1^{t-1}, o_1^t, q_t = j | A, B) \\ = \max_{i=1}^n \underbrace{v_{t-1}(i)}_{\text{prev. Viterbi path prob.}} \underbrace{a_{ij}}_{\text{transition prob.}} \underbrace{b_j(o_t)}_{\text{emission prob.}} \quad (61)$$

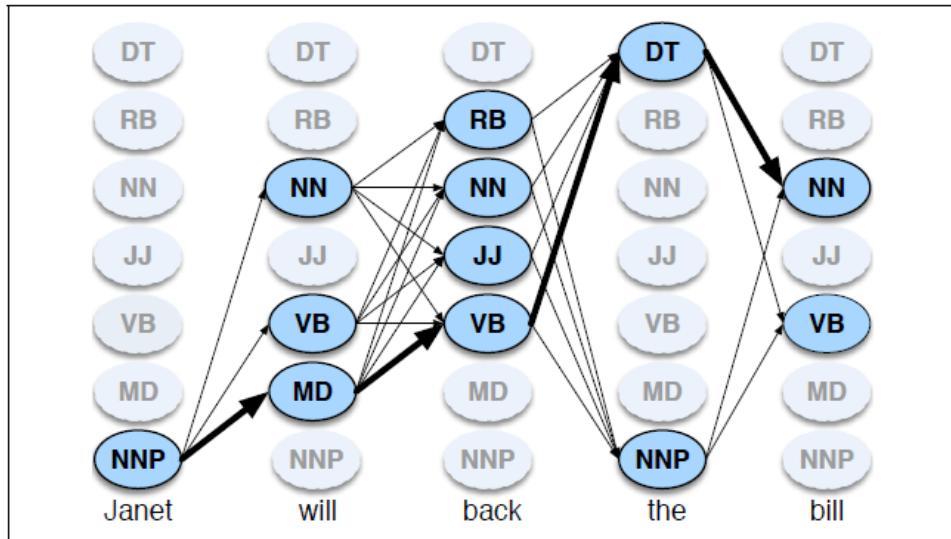


Figure 8.6 A sketch of the lattice for *Janet will back the bill*, showing the possible tags (q_i) for each word and highlighting the path corresponding to the correct tag sequence through the hidden states. States (parts of speech) which have a zero probability of generating a particular word according to the B matrix (such as the probability that a determiner DT will be realized as *Janet*) are greyed out.

This algorithm is slow when the number of states is large ($O(N^2T)$). This can be improved by keeping only the best few hypothesis. One can extend the Viterbi algorithm by using **beam search**, whose approach works as follows:

1. Compute Viterbi scores for each cell.
2. Sort the scores.
3. Keep only the β best-scoring states.

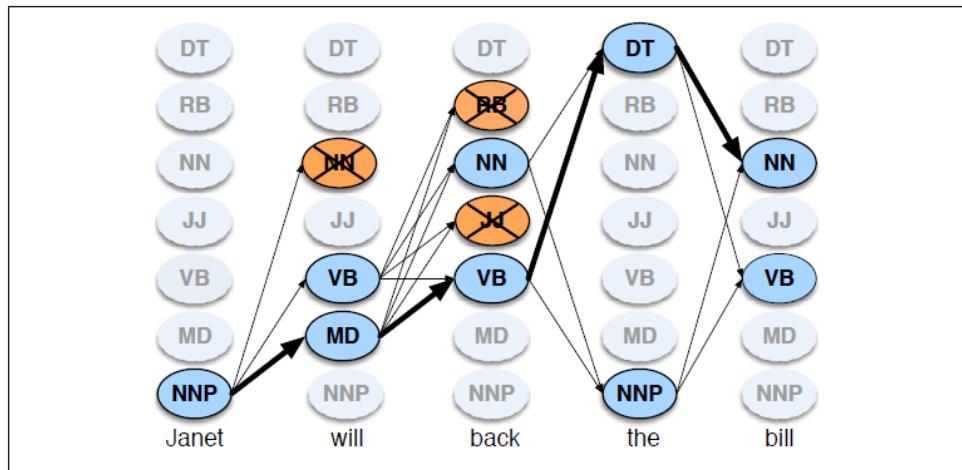
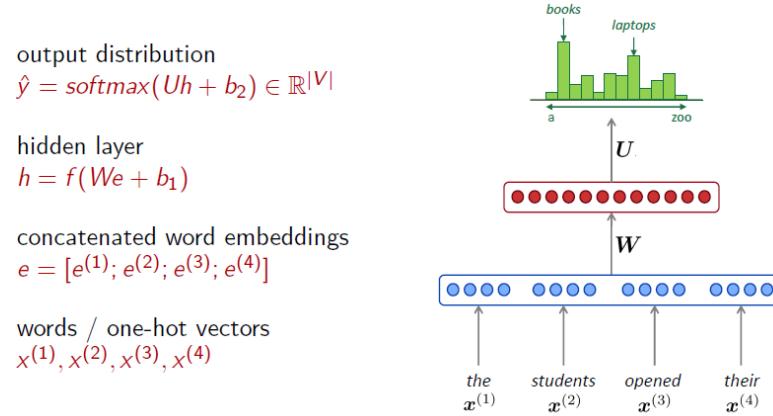


Figure 8.11 A beam search version of Fig. 8.6, showing a beam width of 2. At each time t , all (non-zero) states are computed, but then they are sorted and only the best 2 states are propagated forward and the rest are pruned, shown in orange.

10 Sequence Processing with Recurrent Networks

Language is inherently sequential. Neural language modelling works by applying the following:

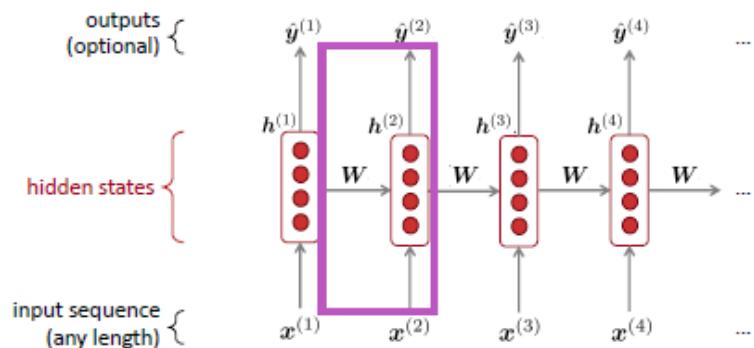


While a fixed-window neural language model is an improvement over N-gram language models, there still exist some problems such as: fixed content size (many language tasks require access to arbitrarily distant information), difficulty in learning systematic patterns.

10.1 Recurrent Neural Networks

... processes appearing in separate windows in consecutive processing steps

The idea behind recurrent neural networks is to apply the same weights W repeatedly:



These networks contain cycles, and the value of a unit depends on the earlier outputs as an input:

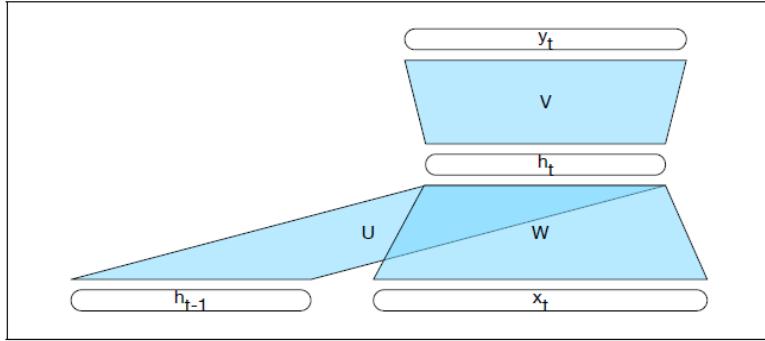


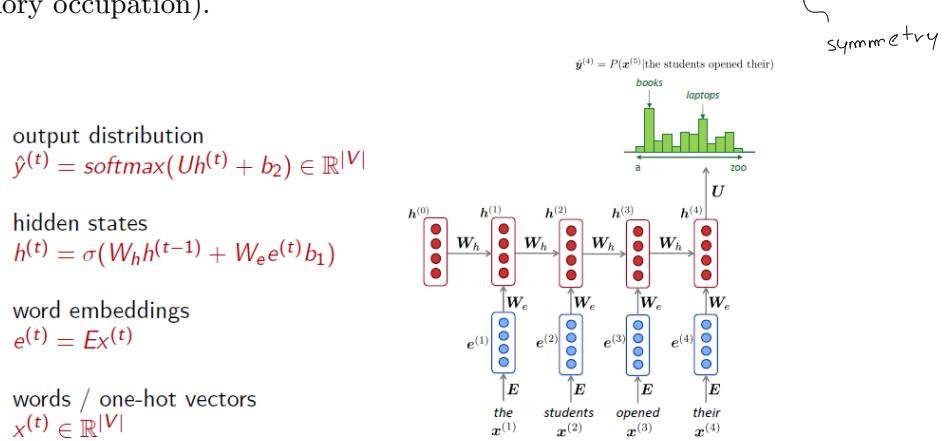
Figure 9.3 Simple recurrent neural network illustrated as a feedforward network.

where:

$$h_t = g(Uh_{t-1} + Wx_t) \quad (63)$$

$$y_t = f(Vh_t) \quad (64)$$

An initial h_0 is given, and f is usually a softmax. In particular, the input vector x_t is multiplied by a weight matrix W and passed through a non-linear activation function g to obtain the activation value for the hidden unit h_t . Such unit is also fed a recurrent link from h_{t-1} . h_{t-1} provides a form of memory or context, encoding earlier processing. Thus, there is no fixed limit on prior context, and the model's size does not increase for longer inputs. Lastly, the same weights are applied on every time step (less memory occupation).



In order to train an RNN language model:

- Get a corpus of text which is a sequence of words $x^{(1)}, \dots, x^{(T)}$.
- Feed into RNN-LM and compute the output distribution $\hat{y}^{(t)}$ for every step t .
- The loss function on step t is the cross-entropy between the predicted probability distributions $\hat{y}^{(t)}$ and the true next word $y^{(t)}$ (one-hot encoding for $x^{(t+1)}$):

41 the input at $t+1$ is
 the true output at t

$$J^{(t)}(\theta) = CE(y^{(t)}, \hat{y}^{(t)}) = - \sum_{w \in V} y_w^{(t)} \log \hat{y}_w^{(t)} = - \log \hat{y}_{x_{t+1}}^{(t)} \quad (65)$$

- Average the loss values to get the overall loss for the entire training set:

$$J(\theta) = \frac{1}{T} \sum_{t=1}^T J^{(t)}(\theta) = \frac{1}{T} \sum_{t=1}^T - \log \hat{y}_{x_{t+1}}^{(t)} \quad (66)$$

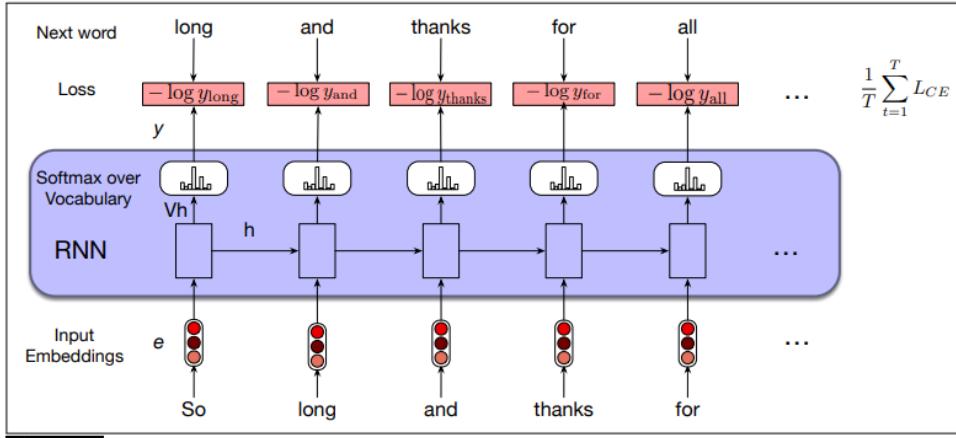


Figure 9.6 Training RNNs as language models.

However, computing the loss and gradients across the entire corpus is usually too expensive. Indeed, the derivative of $J^{(t)}(\theta)$ with respect to the repeated weight matrix W is computed, for every step t , as the sum of the gradients of $J^{(t)}(\theta)$ with respect to W for each step i that precedes t :

$$\frac{\partial J^{(t)}}{\partial W_h} = \sum_{i=1}^t \left. \frac{\partial J^{(t)}}{\partial W_h} \right|_{(i)} \quad (67)$$

Examples of RNN applications are:

- Autoregressive generation (e.g. predicting words). ↘ $x^{(t+1)} \leftarrow \hat{y}^{(t)}$
- Sequence labelling (e.g. POS-tagging). - all output are used
- Sequence classification (e.g. sentiment analysis). - usually, only last hidden state is used
- Text generation (e.g. speech recognition, machine translation, summarization).

10.1.1 Fancier Architectures

Some fancier architectures are the following:

- **Stacked RNN**, where the sequence of outputs from one RNN is used as input to another RNN. This induces representations at different levels of abstraction. However, the training cost rises steeply with the number of levels.

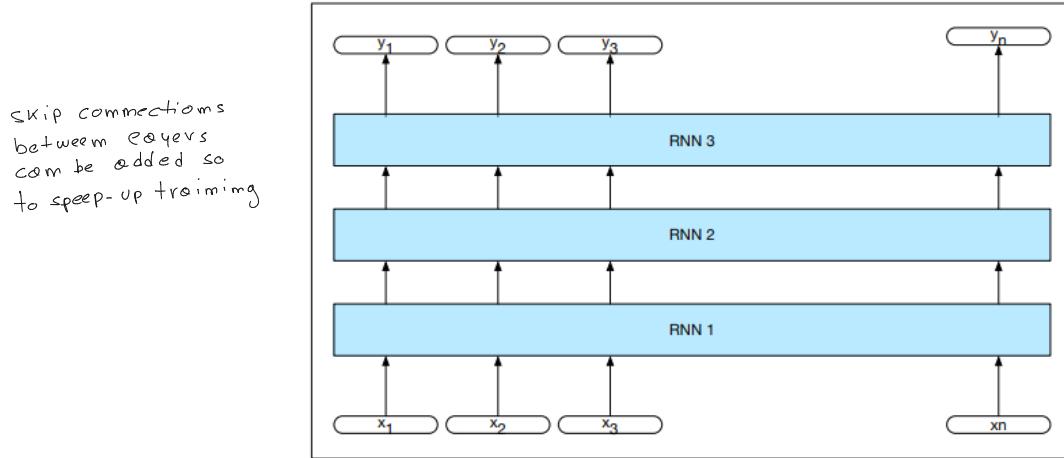


Figure 9.10 Stacked recurrent networks. The output of a lower level serves as the input to higher levels with the output of the last network serving as the final output.

- **Bidirectional RNN**, useful when the entire input sequence is available all at once (usually not for language models). These networks exploit the context to the right of the current input as well. Two independent RNNs, whose outputs are usually concatenated.

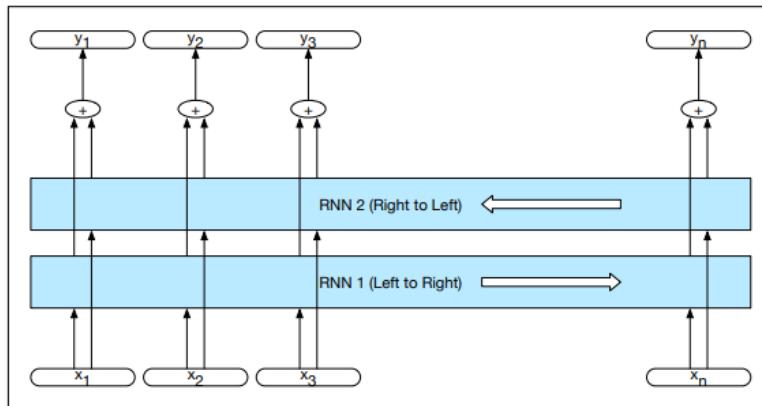


Figure 9.11 A bidirectional RNN. Separate models are trained in the forward and backward directions with the output of each model at each time point concatenated to represent the state of affairs at that point in time. The box wrapped around the forward and backward network emphasizes the modular nature of this architecture.

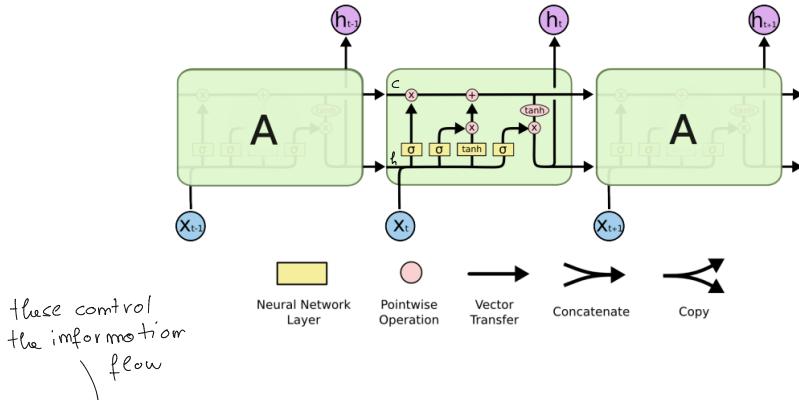
10.2 Exploding and Vanishing Gradient

In general:

- If the gradient becomes too big, then the SGD update step ($\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta)$) becomes too big. A solution to this problem is **gradient clipping** (i.e. scaling down the gradient if its norm exceeds a threshold).
- If the gradient is too small (namely, zero), then the gradient signal from faraway is lost, and the model weights are updated only with respect to near effects, and not to long-term effects. As a solution to the vanishing gradient problem

10.2.1 Long Short-Term Memory

A solution to the vanishing gradient problem are long short-term memory (LSTM) networks. These networks are a type of RNN which, on step t , have an hidden state $h^{(t)}$ and a cell state $c^{(t)}$. Both are vectors of the same length n , and the cell state stores long-term information. In particular, the LSTM can erase, write and read information from the cell state, and the selection of which information to erase/write/read is controlled by three corresponding gates.



The values of the gates are computed based on the current context:

- **Forget gate** controls what is kept and what is forgotten from the previous cell state:

$$f^{(t)} = \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f) \quad (68)$$

- **Input gate** controls what parts of the new cell state content are written to the cell state:

$$i^{(t)} = \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i) \quad (69)$$

- **Output gate** controls what parts of the cell state are output to the hidden state:

$$o^{(t)} = \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o) \quad (70)$$

Moreover:

- **New cell state content** is computed as:

$$\tilde{c}^{(t)} = \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c) \quad (71)$$

- **Cell state**: erase some content from last cell state (forget), and write some new cell content (intput):

$$c^{(t)} = f^{(t)} \cdot c^{(t-1)} + i^{(t)} \cdot \tilde{c}^{(t)} \quad (72)$$

- **Hidden state**: read some content from the cell (output):

$$h^{(t)} = o^{(t)} \cdot \tanh(c^{(t)}) \quad (73)$$

The LSTM architecture makes it easier for the RNN to preserve information over many steps. For example, if the forget state is set to 1 for a cell and the input gate is set to 0, then the information of that cell is preserved indefinitely. By contrast, it is harder for a vanilla RNN to learn a recurrent weight matrix W_h that preserves information in the hidden state.

10.2.2 Gated Recurrent Units

Gated recurrent units (GRU) have been proposed as a simpler alternative to the LSTM. On each step t , one has the input $x^{(t)}$ and the hidden state $h^{(t)}$, and no cell state. Moreover:

- **Update gate** controls what parts of the hidden state are updated and what parts are preserved:

$$u^{(t)} = \sigma(\mathbb{W}_u h^{(t-1)} + U_u x^{(t)} + b_u) \quad (74)$$

- **Reset gate** controls what parts of the previous hidden state are used to compute the new content:

$$r^{(t)} = \sigma(W_r h^{(t-1)} + U_r x^{(t)} + b_r) \quad (75)$$

- **New hidden state content**: the reset gate selects useful parts of the previous hidden state. This is used, alongside the current input, to compute the new hidden content:

$$\tilde{h}^{(t)} = \tanh(W_h(r^{(t)} \cdot h^{(t-1)}) + U_h x^{(t)} + b_h) \quad (76)$$

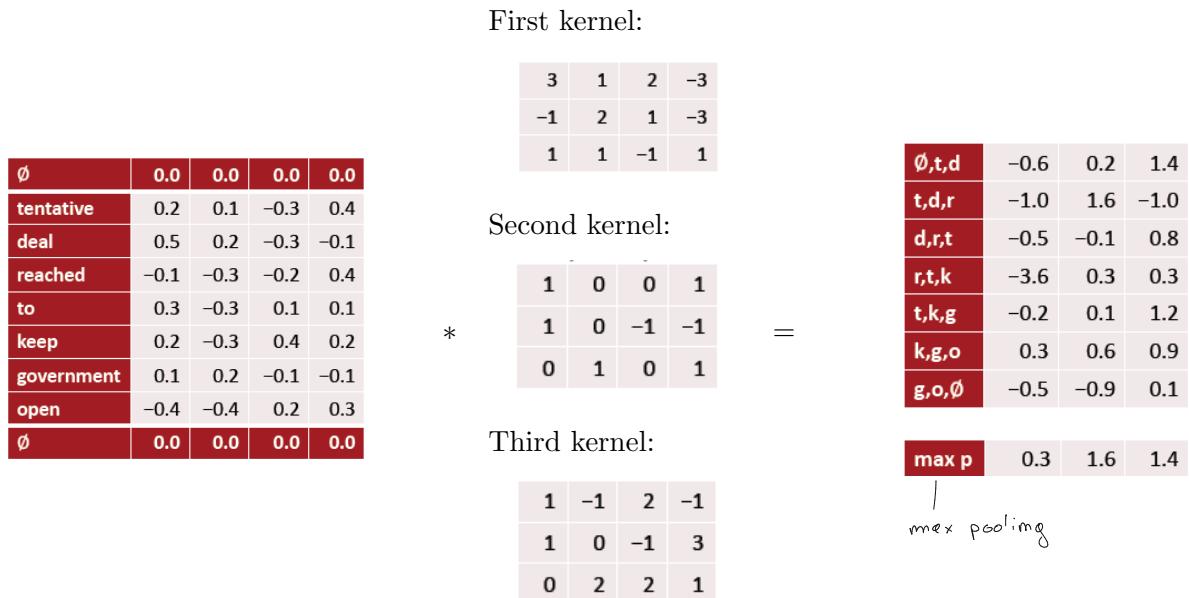
- **Hidden state**: the update gate simultaneously controls what is kept from the previous state, and what is updated in the new hidden state content:

$$h^{(t)} = (1 - u^{(t)}) \cdot h^{(t-1)} + u^{(t)} \cdot \tilde{h}^{(t)} \quad (77)$$

13 CNNs, Transformers and Contextual Word Embeddings

13.1 CNNs for NLP

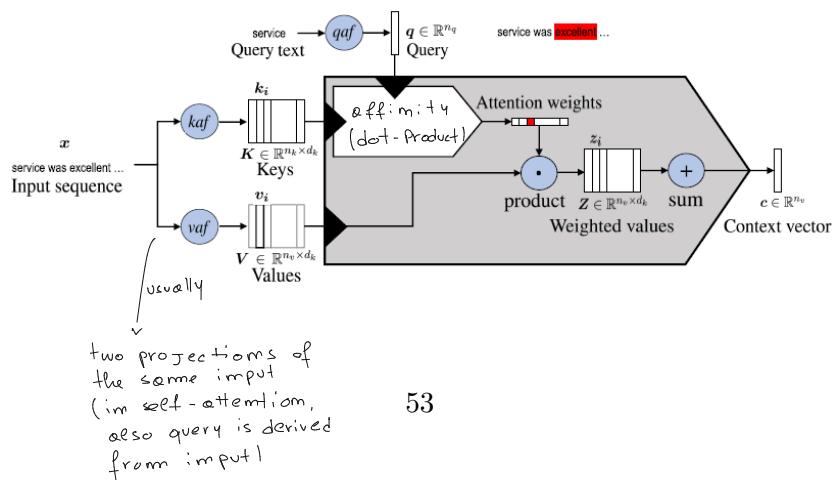
Convolutional neural networks can be applied to text. Indeed, the idea behind 1D convolution is to compute vectors for every possible word subsequence of a certain length. For example, for *tentative deal reached to keep government open*, one computes the vectors for *tentative deal reached*, *deal reached to*, *reached to keep*, *to keep government*, *keep government open*:



In general, CNNs are good for sentence classification but they are somewhat implausible/hard to interpret. Moreover, they are easy to parallelize on GPUs.

13.2 Transformers

Transformers have been introduced to achieve easy parallelization, like CNNs, and maintain the same sequence-to-sequence nature of RNNs. The attention mechanism in transformers uses the concepts of queries, keys and values, which are derived from the input sentence.



The attention mechanism comes in many flavours:

- Often K and V are linear projections of the same input embeddings.
- Q , K and V could also be all projections of the input using different embedding matrices (**self-attention**).
- The combination of Q and K can be a simple dot-product (**dot-product attention**).
- Different projections of Q , K and V can be used to yield multiple outputs (**multi-head attention**).

Transformers were originally introduced in 2017, where the authors proposed to apply a specific version of attention: the scaled dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V \quad \begin{matrix} d_k: \text{dimension of the representation} \\ \text{added to prevent softmax squeezing to zero almost all values} \end{matrix} \quad (82)$$

where a scaling factor is used to prevent extremely small gradients. Moreover, masking can be used to inhibit leftward information flow. The original paper also proposed to use multi-head, scaled dot-product attention. Q , K and V are projected h times with different, learned linear projections:

$$\text{MultiHead}(Q, K, V) = \text{concat}(\text{head}_1, \dots, \text{head}_h)W^Q \quad (83)$$

where:

$$\begin{matrix} \backslash & \text{multiple different projections of the input (parallel computation)} \\ / & \end{matrix} \quad \text{head}_i = \text{Attention}(QW_i^Q, KW_i^K, VW_i^V) \quad (84)$$

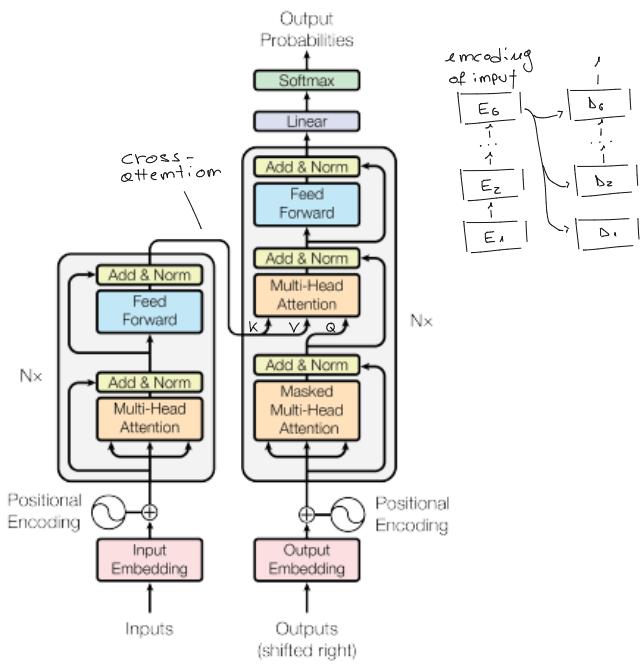
A transformer is an RNN-free, encoder-decoder architecture entirely based on attention:

- The encoder is a stack of N layers, where each layer is composed of two sub-layers:
 - A multi-head self-attention mechanism.
 - A fully connected feed-forward network with ReLU activation, operating on each position independently.

The encoder also uses residual connections, normalization, and positional encoding which makes up for the lack of sequential information processing.

- The decoder is also a stack of N layers, where each layer is composed of three sub-layers:
 - A masked multi-head self-attention layer.
 - A multi-head self-attention layer over the output of the encoder stack.
 - A feed-forward network.

Masking is to ensure that predictions depend only on known outputs. Also the decoder uses positional encoding, residual connections and normalization.



Self-attention directly models relationships between all words in a sentence, regardless of their respective position. This enables learning long-range dependencies.

13.3 Contextual Word Representations

Contextual word vectors are vectors for words in context (word tokens). These vectors depends on arbitrarily long context of nearby words, and not on the context windows. Since RNNs produce context-specific word representations at each position, the idea here is to train a deep bidirectional RNN as language model, then use the context vectors as pre-trained word tokens. An example of such model is BERT.

The encoder part of a transformer can be used to generate an encoding of the input. Differently from GRU and similar, thanks to self-attention, these encodings will be function of the context. Namely, the encoding for the same word can change based on the context. BERT is an example of bidirectional encoder representations from transformers. Pre-training is used to make BERT produce meaningful encodings of the input based on context. Lastly, after pre-training, the model can be fine-tuned for specific applications (e.g. using specific heads).

- Encoders
- Decoders
- Encoders + decoders

For each of these types there exist different applications and pre-training techniques.

Pre-training is based on reconstructing different parts of the input (not only last word) }
These parts of the input have been masked before training (e.g. in BERT 15% masking), }
or changed with other tokens. 55 }

Pre-trained models are large, they are trained on large vocabularies.

There exist benchmarks to evaluate fine-tuned models.

12 Information Extraction and Question Answering

12.1 Information Extraction

Information extraction is about extracting limited kinds of semantic content from text: named entity recognition, relation extraction, event extraction, template filling. In particular:

- **Named entity recognition** is about finding each mention of entities (e.g. people, places, etc.).
For example:

Type	Tag	Sample Categories	Example sentences
People	PER	people, characters	Turing is a giant of computer science.
Organization	ORG	companies, sports teams	The IPCC warned about the cyclone.
Location	LOC	regions, mountains, seas	The Mt. Sanitas loop is in Sunshine Canyon.
Geo-Political	GPE	countries, states, provinces	Palo Alto is raising the fees for parking.
Entity			
Facility	FAC	bridges, buildings, airports	Consider the Golden Gate Bridge.
Vehicles	VEH	planes, trains, automobiles	It was a classic Ford Falcon.

Figure 18.1 A list of generic named entity types with the kinds of entities they refer to.

In particular, given a corpus of text, one has to find proper spans of text that constitute proper names and then classify the type of entity. Named entity recognition can be seen as a sequence labelling task which can be solved using three main approaches: feature-based (HMM, CRF, etc.), neural based (LSTMs), and rule-based.

- **Relation extraction** is about discerning relationships among detected entities. For example:

Relations	Types	Examples
Physical-Located	PER-GPE	He was in Tennessee
Part-Whole-Subsidiary	ORG-ORG	XYZ, the parent company of ABC
Person-Social-Family	PER-PER	Yoko's husband John
Org-AFF-Founder	PER-ORG	Steve Jobs, co-founder of Apple...

Figure 18.10 Semantic relations with examples and the named entity types they involve.

Relations correspond to sets of tuples over elements of a domain. Relation extraction can be solved either using a rule-based approach, or a machine learning approach.

12.2 Question Answering

Question answering focuses on questions that can be answered with simple facts expressed in short texts (factoid questions). There are two paradigms to implement question answering: an IR-based paradigm, which finds relevant documents then uses reading comprehension to read and draw and answer, and a knowledge-based paradigm, which builds a semantic (namely, logic) representation of the query, then queries a database of facts. In particular:

information retrieval

- **Knowledge-based** question answering focuses on semantic parsing to find a mapping between a text string and its logical form. This logic query is translated into a database query.
- **IR-based** question answering focuses on answering a user's question by finding short text segments from a document collection. The main steps that are involved are: question processing, document and passage retrieval, answer extraction.

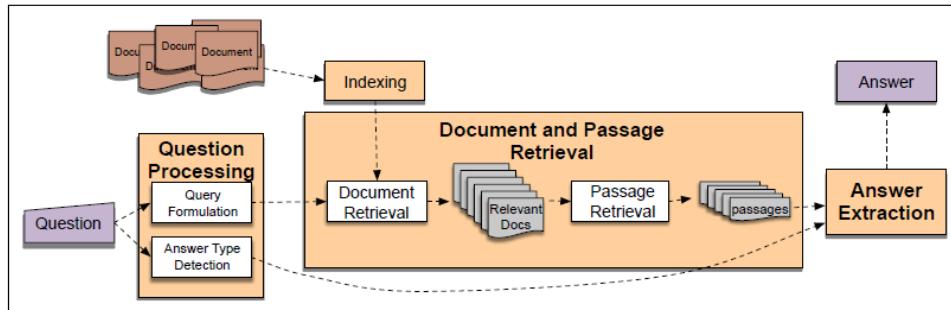


Figure 25.2 IR-based factoid question answering has three stages: question processing, passage retrieval, and answer processing.

- The goal of **question processing** is to extract the query (i.e. a list of keywords for matching documents) and, possibly, the answer type, the focus and the question type.
- **Document and passage retrieval** is about returning a ranked set of documents, and dividing these documents into smaller passages such as sections, paragraphs or sentences.
- **Answer extraction** is a span labelling task: given a text, one would like to find the span of text which constitutes an answer. This type of task can be solved via supervised learning, either feature-based or neural-based. **Feature-based** answer extraction can be carried out using hand-written regular expressions:

Pattern	Question	Answer
<AP> such as <QP>	What is autism?	“, developmental disorders such as autism”
<QP>, a <AP>	What is a caldera?	“the Long Valley caldera, a volcanic crater 19 miles long”

Figure 25.5 Some answer-extraction patterns using the answer phrase (AP) and question phrase (QP) for definition questions (Pasca, 2003).

Neural-based answer extraction, on the other hand, is based on neural networks. This approach works by producing good word embeddings and sentence encodings enabling to select the correct passage spans.

11 Machine Translation, Sequence-to-Sequence and Attention

11.1 Machine Translation

Machine translation is the task of translating a sentence x from one language (the source language) to a sentence y in another language (the target language). Machine translation was first introduced in the early 1950s, using rule-based systems using bilingual dictionaries. Then, in the 1990-2010s period, machine translation was carried out using statistical machine translation, while now neural machine translation is what is mostly used. In particular:

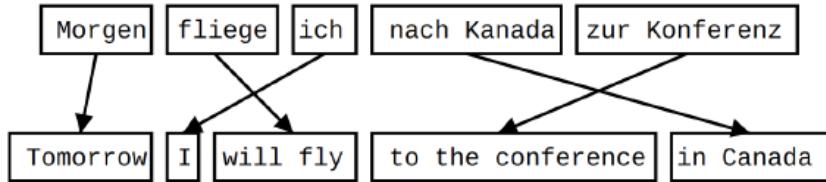
- **Statistical machine translation** is based on learning a probabilistic model from the data: learning which is the best y sentence given the x sentence ($\text{argmax}_y P(y|x)$). This approach uses Bayes rule to split such probability into two components to be learnt separately:
 - Translation model: fidelity, learnt from parallel data.
 - Language model: fluency, learnt from monolingual data.

$$\underset{y}{\text{argmax}} \quad \underbrace{P(x|y)}_{\text{translation model}} \quad \underbrace{P(y)}_{\text{language model}} \quad (78)$$

In particular, the translation model $P(x|y)$ can be learned by breaking it down even further:

$$P(x, a|y) \quad (79)$$

where a is the **alignment**, i.e. the word-level correspondence between source and target:

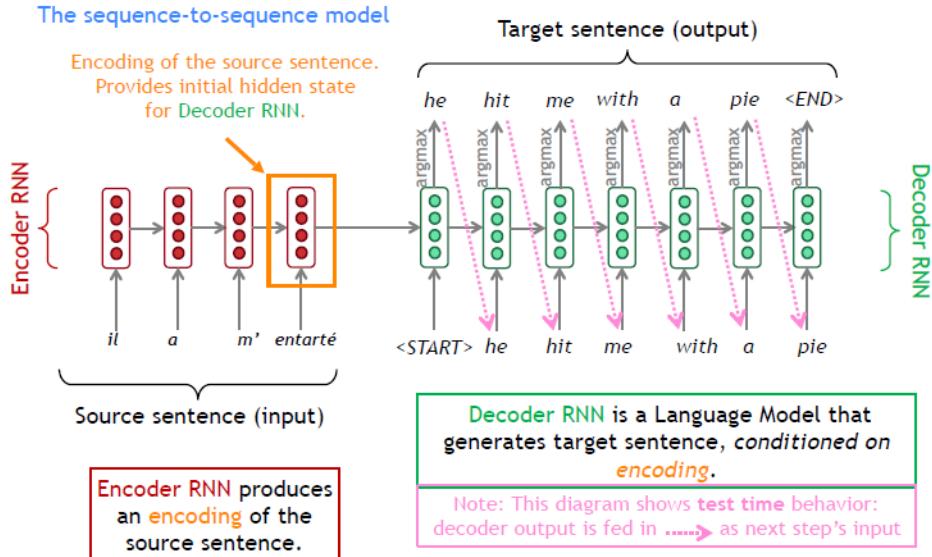


In general, differences in structural/functional features between languages lead to complicated alignments. Moreover, alignment can be one-to-many, and some words do not have a counterpart. $P(x, a|y)$ is learned as a combination of many factors, including:

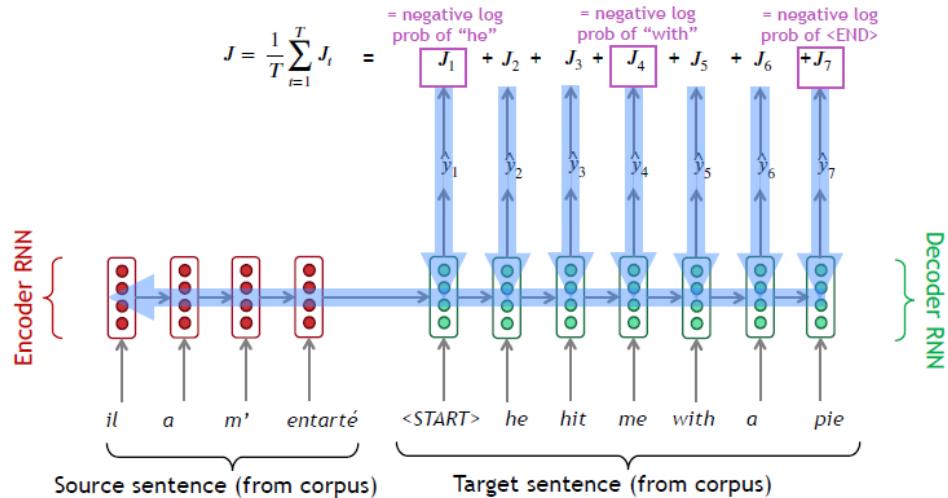
- The probability of particular words aligning.
- The probability of particular words having particular fertility (i.e. number of corresponding words)

Lastly, alignments a are latent variables, thus they are not explicitly specified in the data, and require the use of special learning algorithms for learning the parameters of the distributions with such latent variables.

- **Neural machine translation** is based on doing machine translation using a single neural network which is called sequence-to-sequence or seq2seq. This approach involves two RNNs: one encoder and one decoder.



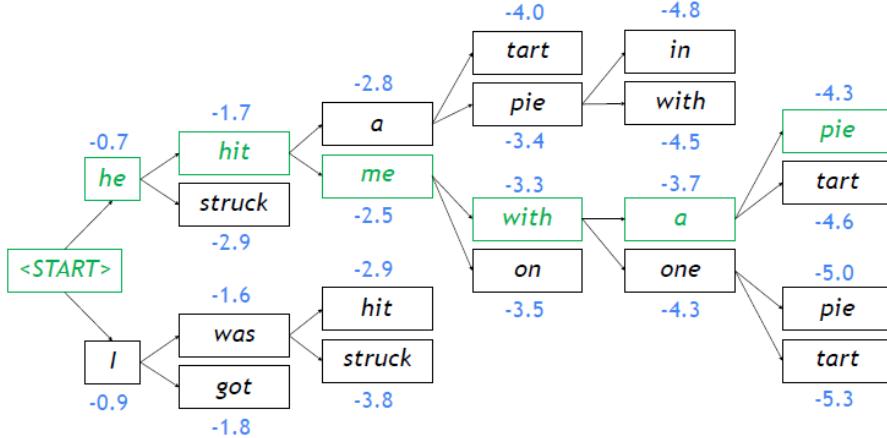
In order to train such systems, one can compute the total loss as:



In the aforementioned case, decoding is applied in a greedy way: the target sentence is generated by taking the argmax on each step of the decoder. However, this could generate a non-optimal translation. To solve this, one could use **beam search decoding**: on each step of the decoder,

one keeps track of the k most probable partial translations (i.e. hypothesis). These hypothesis are given a score of:

$$\text{score}(y_1, \dots, y_t) = \log P_{LM}(y_1, \dots, y_t | x) \quad (80)$$



Different hypothesis may produce the <END> tokens on different steps. An hypothesis that produced <END> is marked as complete. Complete hypothesis are placed aside and other hypothesis are explored, using beam search, until n hypothesis, or a maximum predefined length, are reached. Finally, the best hypothesis are selected by computing $\text{score}(y_1, \dots, y_t)/t$.

11.2 Evaluation Metrics

Machine translation can be evaluated in different ways:

- **Bilingual Evaluation Understudy (BLEU)** compares the candidate machine-written translation to one or several reference human-written translation(s). The proposed similarity score is based on:

- N-gram precision score (usually for 1, 2, 3 and 4-grams): the fraction of n-grams in the candidate translation that also appear in the reference translation.
- Count at most one match for each n-gram in the reference translation.
- Add a penalty for too-short system translations.

BLEU is useful but imperfect. There are many valid ways to translate a sentence, so a good translation can get a poor BLEU score because it has low n-gram overlap with the human translation.

- **Character F-score (chrF)** is a more recent alternative to BLEU. It is based on character n-grams:

- $chrP$ is the percentage of char 1-grams, 2-grams, …, k -grams in the hypothesis that occur in the reference, averaged.
- $chrR$ is the percentage of char 1-grams, 2-grams, …, k -grams in the reference that occur in the hypothesis, averaged.

$$chrF\beta = (1 + \beta^2) \frac{chrP \cdot chrP}{\beta^2 \cdot chrP + chrR} \quad (81)$$

- **Embedding similarity**, which is more robust against synonyms or paraphrases.

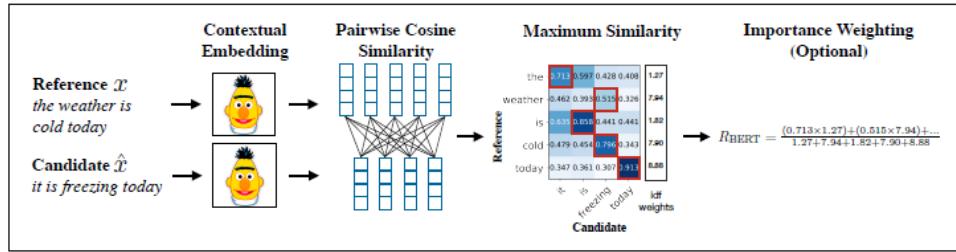
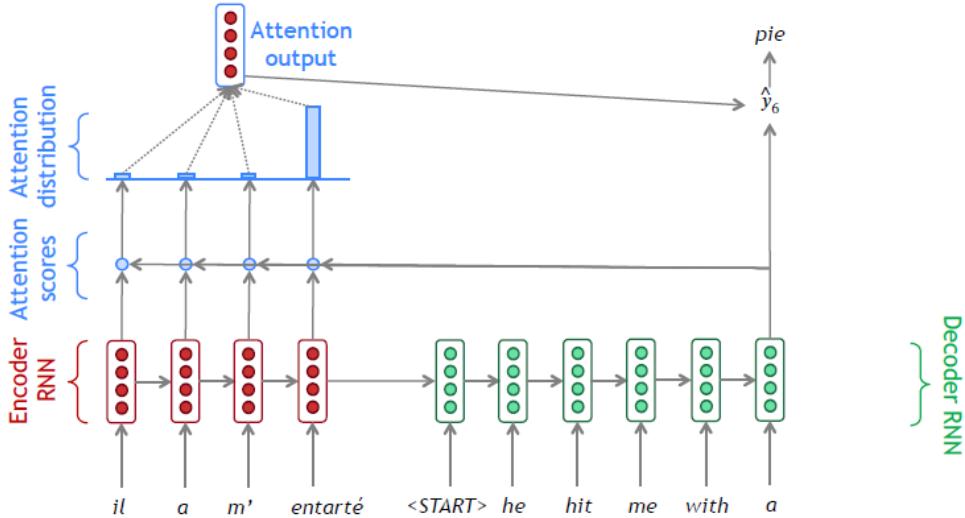


Figure 10.18 The computation of BERTSCORE recall from reference x and candidate \hat{x} , from Figure 1 in Zhang et al. (2020). This version shows an extended version of the metric in which tokens are also weighted by their idf values.

11.3 Attention

On each step of the decoder, one could use a direct connection to the encoder to focus on a particular part of the source sequence:



In particular:

- The encoder hidden states are $h_1 \dots h_N \in \mathbb{R}^h$.
- The decoder hidden state at step t is s_t .
- The **attention scores** at step t are computed as $e^t = [s_t^T h_1, \dots, s_t^T h_N] \in \mathbb{R}^N$.
- The **attention distribution** α^t used to take a weighted sum of the encoder hidden states is computed as $\alpha^t = \text{softmax}(e^t) \in \mathbb{R}^N$.
- The **attention output**, which mostly contains information from the hidden states the received high attention, is computed as $a_t = \sum_{i=1}^N \alpha_i^t h_i \in \mathbb{R}^h$.
- The attention output and the decoder hidden state are concatenated $[a_t; s_t] \in \mathbb{R}^{2h}$.

14 Dialogue Systems and Chatbots

14.1 Dialogue Systems

Some Characteristics of human conversations are the following:

- Turn, namely each single contribution to the dialogue. A system has to know when to start/stop talking, and has to perform endpoint detection.
- Speech act, namely the action implied by a dialogue utterance.
- Grounding, namely establishing common ground.
- Dialogue structure, namely the structure followed by dialogue turns.
- Initiative, namely who controls the conversation.
- Implicature, namely a licensed inference (i.e. what inference is the hearer expected to draw? Usually, more information is communicated than present in uttered words).

In general, there exist two main types of dialogue systems: **chatbots**, which can carry on extended conversations with the goal of mimicking the unstructured conversations characteristics of the informal human-human interaction (e.g. ELIZA), and **task-oriented dialogue agents**, which use conversation with users to help complete tasks (e.g. Siri, Alexa, etc.).

14.2 Chatbots

The most important chatbot dialogue system in the field's history is ELIZA. ELIZA is designed to simulate a Rogerian psychologist (i.e. the system can assume the pose of knowing almost nothing about the world). This chatbot is a rule-based system. There exist also corpus-based chatbots, which are usually built using IR and machine-learned sequence transduction:

- information retrieval → $\xrightarrow{\text{encoder-decoder architectures}}$
- **IR-based chatbots** respond to user's turn (query q) by repeating some appropriate turn (response r) from a corpus of natural text. In particular, there are two methods to do so:

- Return the response to the most similar turn:

$$r = \text{response} \left(\underset{t \in C}{\operatorname{argmax}} \frac{q^T t}{\|q\| \|t\|} \right) \quad (85)$$

- Return the most similar turn:

$$r = \underset{t \in C}{\operatorname{argmax}} \frac{q^T t}{\|q\| \|t\|} \quad (86)$$

Returning the most similar turn works better in practice.

14.3 Task-Oriented Dialogue Agents

These agents are also called **frame-based dialogue systems**. A frame is a knowledge structure representing the kinds of intentions the system can extract from user sentences. A frame consists of a collection of slots, representing what the system needs to know. Each slot can take a set of possible values of a particular semantic type. A set of such frames is called **domain ontology**:

Slot	Type	Question Template
ORIGIN CITY	city	“From what city are you leaving?”
DESTINATION CITY	city	“Where are you going?”
DEPARTURE TIME	time	“When would you like to leave?”
DEPARTURE DATE	date	“What day would you like to leave?”
ARRIVAL TIME	time	“When do you want to arrive?”
ARRIVAL DATE	date	“What day would you like to arrive?”

Figure 26.9 A frame in a frame-based dialogue system, showing the type of each slot and a question used to fill the slot.

The system’s goal is to fill the slots in the frame, then perform an action which can be relevant for the user. In particular, the system asks questions to the user, using question templates. If the user’s response fills multiple slots, the system fills all relevant slots. Once enough information have been collected, the system can perform a necessary action and return the result. A common approach to slot-filling is rule-based, where the system consists of a large hand-designed semantic grammar (i.e. CFG where the left-hand side of each rule corresponds to semantic entities).

SHOW	→ show me i want can i see ...
DEPART_TIME_RANGE	→ (after around before) HOUR morning afternoon evening
HOUR	→ one two three four... twelve (AMPM)
FLIGHTS	→ (a) flight flights
AMPM	→ am pm
ORIGIN	→ from CITY
DESTINATION	→ to CITY
CITY	→ Boston San Francisco Denver Washington

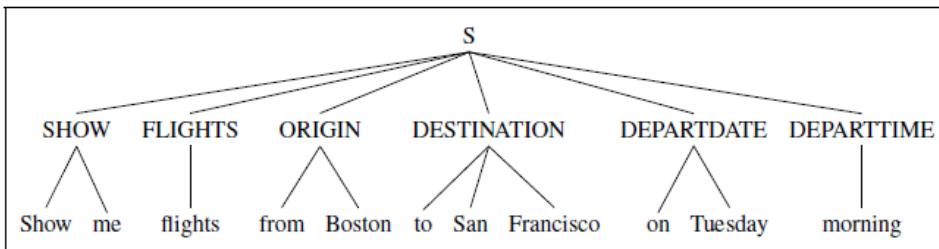


Figure 26.10 A semantic grammar parse for a user sentence, using slot names as the internal parse tree nodes.

More advanced architectures have also other components: natural language understanding, which helps extracting slot fillers using machine learning, dialogue state tracker, which maintains the current state

of dialogue, dialogue policy, which decides what to do or say next, and natural language generation, which is conditioned on context to produce more natural turns.

14.4 Debating Systems

The general field is argument miming. These systems are based on debate principles (claims, evidences, arguments).