

Linear Models for Statistical Natural Language Processing

Jacob Eisenstein

October 4, 2015

Contents

Contents	1
I Words, bags of words, and features	11
1 Linear classification and features	13
1.1 Review of basic probability	15
1.2 Naïve Bayes	20
1.3 Recap	28
2 Discriminative learning	31
2.1 Features	31
2.2 Perceptron	32
2.3 Loss functions and large-margin classification	33
2.4 Logistic regression	38
2.5 Summary of learning algorithms	45
3 Linguistic applications of classification	49
3.1 Sentiment and opinion analysis	49
3.2 Word sense disambiguation	52
3.3 Other applications	56
4 Learning without supervision	57
4.1 K-means clustering	58
4.2 The Expectation-Maximization (EM) Algorithm	59
4.3 Applications of EM	65
4.4 Other approaches to learning with latent variables*	68
5 Language models	71

5.1	N-gram language models	74
5.2	Evaluating language models	76
5.3	Smoothing and discounting	80
5.4	Other Types of N-gram Language Models	84
5.5	Other details	88
II	Sequences and trees	91
6	Morphology	93
6.1	Types of morphemes	97
6.2	Types of morphology	98
6.3	Computing and morphology	105
7	Finite-state automata	107
7.1	Automata and languages	108
7.2	Weighted Finite State Automata	114
7.3	Semirings	118
7.4	Finite state transducers	121
7.5	Weighted FSTs	123
7.6	Applications of finite state composition	125
7.7	Discriminative structure prediction	127
8	Part-of-speech tagging	129
8.1	Details about parts-of-speech	132
8.2	Part of speech tagging	136
9	Sequence labeling	141
9.1	Hidden Markov Models	142
9.2	Algorithms for sequence labeling	144
9.3	Learning models of sequence labeling	151
9.4	Conditional random fields	155
9.5	Unsupervised sequence labeling*	164
10	Context-free grammars	167
10.1	Is English a regular language?	167
10.2	Context-Free Languages	169
10.3	Constituents	171
10.4	A simple grammar of English	174

10.5 Grammar design	177
10.6 Grammar equivalence and normal form	177
11 CFG Parsing	179
11.1 CKY parsing	179
11.2 Ambiguity in parsing	181
11.3 Statistical parsing with PCFGs	185
11.4 Algorithms for PCFG Parsing	188
11.5 Parser evaluation	189
11.6 Improving PCFG parsing	190
11.7 Modern constituent parsing	205
12 Dependency Parsing	211
12.1 Algorithms	215
12.2 Higher-order dependency parsing	221
12.3 Learning dependency parsers	222
12.4 Applications	223
III Meaning	225
13 Logical semantics	227
14 Shallow semantics	229
14.1 Predicates and arguments	229
14.2 Resources for shallow semantics	233
14.3 FrameNet	234
14.4 Semantic Role Labeling	236
14.5 Abstract Meaning Representation	237
15 Distributional semantics	239
15.1 The distributional hypothesis	239
15.2 Local context	241
15.3 Syntactic context	244
15.4 Latent semantic analysis	245
15.5 Word vectors and neural word embeddings	246
16 Discourse	251
16.1 Discourse relations in the Penn Discourse Treebank	251

16.2 Rhetorical Structure Theory	251
16.3 Centering	251
16.4 Lexical cohesion and text segmentation	252
16.5 Dialogue	252
17 Anaphora and Coreference Resolution	253
17.1 Forms of referring expressions	254
17.2 Pronouns and reference	257
17.3 Resolving ambiguous references	259
17.4 Coreference resolution	262
17.5 Coreference evaluation	265
17.6 Multidocument coreference resolution	265
 IV Applications	 267
18 Information extraction	269
18.1 Entities	269
18.2 Relations	269
18.3 Events and processes	269
18.4 Facts, beliefs, and hypotheticals	269
19 Machine translation	271
19.1 The noisy channel model	271
19.2 Translation modeling	273
19.3 Example for IBM Model 1	276
19.4 Syntactic MT	276
19.5 Algorithms for SCFGs	280
 V Learning	 283
20 Semi-supervised learning	285
20.1 Learning with less annotation effort	286
20.2 Why would unlabeled data help?	287
20.3 Domain adaptation	295
20.4 Other learning settings	298
21 Beyond linear models	299

<i>CONTENTS</i>	5
21.1 Representation learning	299
21.2 Convolutional neural networks	299
21.3 Recursive neural networks	299
21.4 Encoder-decoder models	299
Bibliography	301

Preface

This text is built from the notes that I use for teaching Georgia Tech’s undergraduate and graduate courses on natural language processing, CS 4650 and 7650. The focus is on what I view as a core subset of the field of natural language processing, unified by the concepts of linear models and structure prediction. This includes approaches to document classification, word sense disambiguation, sequence labeling (part-of-speech tagging and named entity recognition), parsing, coreference resolution, relation extraction, discourse analysis, and, to a limited degree, language modeling and machine translation. The title was inspired by Fernando Pereira’s EMNLP 2008 keynote, “Are linear models right for language.”¹ The notes are influenced by several other good resources (e.g., Manning and Schütze, 1999; Jurafsky and Martin, 2009; Smith, 2011; Figueiredo et al., 2013; Collins, 2013), but for various reasons I wanted to create something of my own.

The text assumes familiarity with basic linear algebra, and with calculus through Lagrange multipliers. It includes a refresher on probability, but some previous exposure would be helpful. An introductory course on the analysis of algorithms is also assumed; in particular, the reader should be familiar with asymptotic analysis of the speed and memory costs of algorithms, and should have seen dynamic programming. No prior background in machine learning or linguistics is assumed, and even students with background in machine learning should be sure to read the introductory chapters, since the notation used in natural language processing is different from typical presentations of machine learning classifiers, due to the heavy emphasis on structure prediction in applications of machine learning to language. Throughout the book, advanced material is marked with an asterisk, and can be safely skipped.

-Jacob Eisenstein, October 4, 2015

¹You can see a version of this talk — not the one I saw — online at vimeo.com/30676245

Notation

w_n	word token at position n
\mathbf{x}_i	a vector of feature counts for instance i , often word counts
N	number of training instances
V	number of words in vocabulary
$\boldsymbol{\theta}$	a vector of weights
y_i	the label for instance i
\mathbf{y}	vector of labels across all instances
\mathcal{Y}	set of all possible labels
K	number of possible labels $K = \# \mathcal{Y} $
$\mathbf{f}(\mathbf{x}_i, y_i)$	feature vector for instance i with label y_i
$P(A)$	probability of event A
$p_B(b)$	the marginal probability of random variable B taking value b
M	length of a sequence (of words or tags)
$\mathcal{T}(\mathbf{w})$	the set of possible tag sequences for the word sequence \mathbf{w}
$\langle \text{START} \rangle$	the start symbol
$\langle \text{STOP} \rangle$	the stop symbol
λ	the amount of regularization

Part I

Words, bags of words, and features

Chapter 1

Linear classification and features

Suppose you want to build a spam detector, in which each document is classified as “spam” or “ham.” How would you do it, using only the text in the email?

One solution is to represent document i as a column vector of word counts: $\mathbf{x}_i = [0 \ 1 \ 1 \ 0 \ 0 \ 2 \ 0 \ 1 \ 13 \ 0 \ \dots]^\top$, where $x_{i,j}$ is the count of word j in document i . Suppose the size of the vocabulary is V , so that the length of \mathbf{x}_i is also V . The object \mathbf{x}_i is a vector, but colloquially we call it a **bag of words**, because it includes only information about the count of each word, and not the order in which they appear.

We’ve thrown out grammar, sentence boundaries, paragraphs — everything but the words! But this could still work. If you see the word *free*, is it spam or ham? How about *Bayesian*? One approach would be to define a “spamminess” score for every word in the dictionary, and then just add them up. These scores are called **weights**, written θ , and we’ll spend a lot of time talking about where they come from.

But for now, let’s generalize: suppose we want to build a multi-way classifier to distinguish stories about sports, celebrities, music, and business. Each label is an element y_i in a set of K possible labels \mathcal{Y} . Our goal is to predict a label \hat{y}_i , given the bag of words \mathbf{x}_i , using the weights θ . We’ll do this using a vector inner product between the weights θ and a **feature vector** $\mathbf{f}(\mathbf{x}_i, y_i)$. As the notation suggests, the feature vector is constructed by combining \mathbf{x}_i and y_i . For example, feature j might be,

$$f_j(\mathbf{x}_i, y_i) = \begin{cases} 1, & \text{if } (\text{free} \in \mathbf{x}_i) \wedge (y_i = \text{SPAM}) \\ 0, & \text{otherwise} \end{cases} \quad (1.1)$$

For any pair $\langle \mathbf{x}_i, y_i \rangle$, we then define $\mathbf{f}(\mathbf{x}_i, y_i)$ as,

$$\mathbf{f}(\mathbf{x}, Y = 0) = [\mathbf{x}^\top \mathbf{0}_{V(K-1)}^\top]^\top \quad (1.2)$$

$$\mathbf{f}(\mathbf{x}, Y = 1) = [\mathbf{0}_V^\top \mathbf{x}^\top \mathbf{0}_{V(K-2)}^\top]^\top \quad (1.3)$$

$$\mathbf{f}(\mathbf{x}, Y = 2) = [\mathbf{0}_{2V}^\top \mathbf{x}^\top \mathbf{0}_{V(K-3)}^\top]^\top \quad (1.4)$$

...

$$\mathbf{f}(\mathbf{x}, Y = K) = [\mathbf{0}_{V(K-1)}^\top \mathbf{x}^\top]^\top, \quad (1.5)$$

where $\mathbf{0}_{VK}$ is a column vector of $V \times K$ zeros. This arrangement is shown in Figure 1.1. This notation may seem like a strange choice, but in fact it helps to keep things simple. Given a vector of weights, $\boldsymbol{\theta} \in \mathbb{R}^{V \times K}$, we can now compute the inner product $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$. This inner product gives a scalar measure of the score for label y , given observations \mathbf{x} . For any document \mathbf{x}_i , we predict the label \hat{y} as

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y) \quad (1.6)$$

This inner product is the fundamental equation for linear classification, and it is the reason we prefer the feature function notation $\mathbf{f}(\mathbf{x}, y)$. The notation gives a clean separation between the **data** $\mathbf{f}(\mathbf{x}, y)$, and the **parameters**, which are expressed by the single vector of weights, $\boldsymbol{\theta}$. As we will see in later chapters, it generalizes nicely to **structured output spaces**, in which the space of labels \mathcal{Y} is very large, and we want to model shared substructure between labels.

Often we'll add an **offset** feature at the end of \mathbf{x} , which is always 1; we then have to also add an extra zero to each of the zero vectors. This gives the entire feature vector $\mathbf{f}(\mathbf{x}, y)$ a length of $(V+1) \times K$. The weight associated with this offset feature can be thought of as a "bias" for each label. For example, if we expect most documents to be spam, then the weight for the offset feature for $Y = \text{spam}$ should be larger than the weight for the offset feature for $Y = \text{ham}$.

Returning to the weights $\boldsymbol{\theta}$ — where do they come from? As already suggested, we could just set the weights by hand. If we wanted to distinguish, say, English from Spanish, we could just use English and Spanish dictionaries, and set the weight to 1 for each word that appears in the associated dictionary. For example,

$$\begin{array}{ll} \theta_{\text{english}, \text{bicycle}} = 1 & \theta_{\text{spanish}, \text{bicycle}} = 0 \\ \theta_{\text{english}, \text{bicicleta}} = 0 & \theta_{\text{spanish}, \text{bicicleta}} = 1 \\ \theta_{\text{english}, \text{con}} = 1 & \theta_{\text{spanish}, \text{con}} = 1 \\ \theta_{\text{english}, \text{ordinateur}} = 0 & \theta_{\text{spanish}, \text{ordinateur}} = 0 \end{array}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

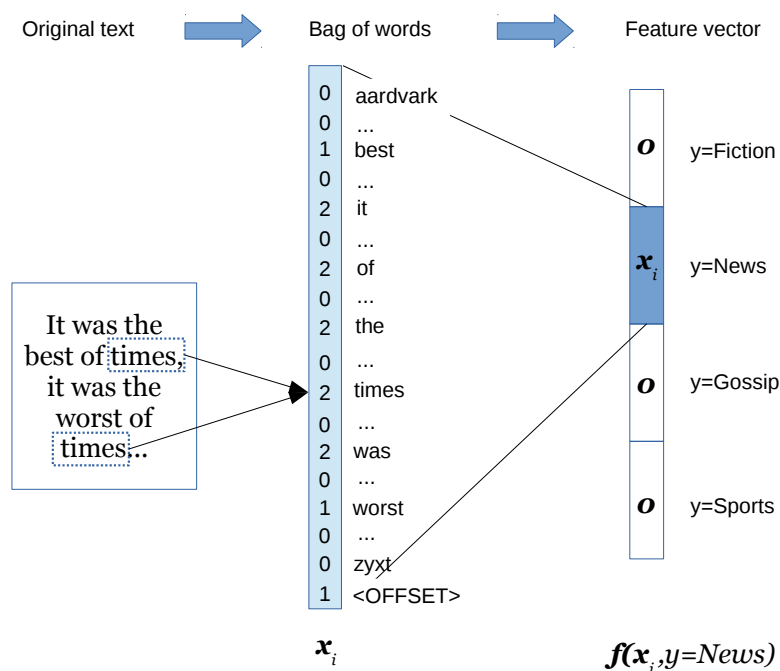


Figure 1.1: The bag-of-words and feature vector representations, for a hypothetical text classification task.

Similarly, if we want to distinguish positive and negative sentiment, we could use positive and negative *sentiment lexicons*, which are defined by expert psychologists (Tausczik and Pennebaker, 2010). You'll try this in Project 1.

But it is usually not easy to set the weights by hand. Instead, we will learn them from data. For example, suppose that an email user has manually labeled thousands of messages as "spam" or "not spam"; or a newspaper may label its own articles as "business" or "fashion." Such **instance labels** are a typical form of labeled data that we will encounter in NLP. In **supervised machine learning**, we use instance labels to automatically set the weights for a classifier. An important tool for this is probability.

1.1 Review of basic probability

This section is inspired by and partially borrowed from Manning and Schütze (1999). If you feel very confident in your understanding of probability, feel free to skim ahead to Section 1.2, where we return to text classification.

- **Formally:** When we write $P(\cdot)$, this denotes a function $P : \mathcal{F} \rightarrow [0, 1]$ from an **event space** \mathcal{F} to a **probability**. A probability is a real number between zero and one, with zero representing impossibility and one representing certainty.
- We think about the event space \mathcal{F} as a set, with any element $A \in \mathcal{F}$ referred to as an **event**. We write \emptyset to indicate the impossible event, $P(\emptyset) = 0$, and Ω to indicate the certain event, $P(\Omega) = 1$.
- If $A_i \in \mathcal{F}$ and $A_j \in \mathcal{F}$ and $A_i \cap A_j = \emptyset$, then A_i and A_j are **disjoint** events. Consider rolling a die, with A_i being the event of rolling 1, and A_j being the event of rolling 2; these are disjoint events, $A_i \cap A_j = \emptyset$. On other hand, if A_i is the event of there being an earthquake, and A_j is the event of there being a hurricane, $A_i \cap A_j \neq \emptyset$, because it is possible to have both an earthquake and a hurricane.
- The probabilities of disjoint event sets are additive:

$$A_i \cap A_j = \emptyset \Rightarrow P(A_i \cup A_j) = P(A_i) + P(A_j). \quad (1.7)$$

This is a restatement of the Third Axiom of probability, which generalizes to any countable sequence of disjoint event sets.

- As an example, you might ask what is the probability of two heads on three flips of a fair coin. There are eight possible series of three flips HHH, HHT, \dots , and each is an equally likely event, with probability $\frac{1}{8}$. Of these events, three meet the criterion of having two heads: HHT, HTH, THH . These events are all mutually exclusive; in other words, each pair of events is disjoint. So the probability is $\frac{1}{8} + \frac{1}{8} + \frac{1}{8} = \frac{3}{8}$.
- More generally, $P(A_i \cup A_j) = P(A_i) + P(A_j) - P(A_i \cap A_j)$. This can be derived from the Third Axiom of probability, mentioned above.

$$P(A_i \cup A_j) = P(A_i) + P(A_j - (A_i \cap A_j)) \quad (1.8)$$

$$P(A_j) = P(A_j - (A_i \cap A_j)) + P(A_i \cap A_j) \quad (1.9)$$

$$P(A_j - (A_i \cap A_j)) = P(A_j) - P(A_i \cap A_j) \quad (1.10)$$

$$P(A_i \cup A_j) = P(A_i) + P(A_j) - P(A_i \cap A_j) \quad (1.11)$$

- If the probability $P(A \cap B) = P(A)P(B)$, then the events A and B are *independent*, written $A \perp B$.

Conditional probability and Bayes' rule

A conditional probability is an expression like $P(A | B)$, where we are interested in the probability of A conditioned on B happening: for example, the probability of a randomly selected person answering the phone by saying *hello*, conditioned on that person being a speaker of English. We define conditional probability as the ratio,

$$P(A | B) = \frac{P(A \cap B)}{P(B)} \quad (1.12)$$

The **chain rule** states that $P(A \cap B) = P(A | B)P(B)$, which is just a simple rearrangement of terms from Equation 1.12. We can apply the chain rule multiple times:

$$\begin{aligned} P(A \cap B \cap C) &= P(A | B \cap C)P(B \cap C) \\ &= P(A | B \cap C)P(B | C)P(C) \end{aligned}$$

Bayes' rule (sometimes called Bayes' law or Bayes' theorem) follows from the chain rule:

$$P(A | B) = \frac{P(A \cap B)}{P(B)} = \frac{P(B | A)P(A)}{P(B)} \quad (1.13)$$

The terms in Bayes rule have specialized names, which we will occasionally use:

- $P(A)$ is the **prior**, since it is the probability of event A without knowledge about whether B happens or not.
- $P(B | A)$ is the **likelihood**, the probability of event B given that event A has occurred.
- $P(A | B)$ is the **posterior**, since it is the probability of event A with knowledge that B has occurred.

Often we want the maximum a posteriori (MAP) estimate,

$$\begin{aligned} \hat{A} &= \arg \max_A P(A | B) \\ &= \arg \max_A P(B | A)P(A)/P(B) \\ &\propto \arg \max_A P(B | A)P(A). \end{aligned}$$

- We don't need to normalize the probability because $P(B)$ is the same for all values of A .
- If we do need to compute the normalized probability $P(A | B)$, we can compute $P(B)$ by summing over $P(B \cap A) + P(B \cap \bar{A})$, where \bar{A} is the **complement** of A . The complement is defined such that $A \cap \bar{A} = \emptyset$ and $A \cup \bar{A} = \Omega$, so that $P(A \cap \bar{A}) = 0$ and $P(A \cup \bar{A}) = 1$.
- More generally, if $\bigcup_i A_i = \Omega$ and $\forall_{i,j}, A_i \cap A_j = \emptyset$, then

$$P(B) = \sum_i P(B | A_i)P(A_i). \quad (1.14)$$

Example Manning and Schütze (1999) have a nice example of Bayes' rule (sometimes called Bayes Law) in a linguistic setting. (This same example is usually framed in terms of tests for rare diseases.)

- Suppose one is interested in a rare syntactic construction, such as parasitic gaps, which occurs on average once in 100,000 sentences.
 - (An example of a sentence with a parasitic gap is *Which class did you attend ... without registering for ...?*)
- Lana Linguist has developed a complicated pattern matcher that attempts to identify sentences with parasitic gaps. Its pretty good, but it's not perfect:
 - If a sentence has a parasitic gap, the pattern matcher will find it with probability 0.95. This is the **recall**; the **false negative rate** is defined as one minus the recall.
 - If the sentence doesn't have a parasitic gap, the pattern matcher will wrongly say it does, with probability 0.005. This is the **false positive rate**.
- Suppose the test says that a sentence contains a parasitic gap. What is the probability that this is true?

Solution: Let G be the event of a sentence having a parasitic gap, and T be the event of the test being positive.

$$P(G | T) = \frac{P(G | T)P(T)}{P(G | T)P(T) + P(G | \bar{T})P(\bar{T})} \quad (1.15)$$

$$= \frac{0.95 \times 0.00001}{0.95 \times 0.00001 + 0.005 \times 0.99999} \approx 0.002 \quad (1.16)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

Note that even though the pattern matcher is very accurate, with false positive and false negative rates below 5%, the extreme rarity of this phenomenon means that a positive result from the detector is most likely to be wrong.

If $P(A \cap B \mid C) = P(A \mid C)P(B \mid C)$, then the events A and B are **conditionally independent**, written $A \perp B \mid C$.

Random variables

A random variable takes on a specific value in \mathbb{R}^n , typically with $n = 1$, but not always. Discrete random variables can take values only in some countable subset of \mathbb{R} .

- Recall the coin flip example. The number of heads, H , can be viewed as a discrete random variable, $H \in 0, 1, 2, 3$.
- Each possible value is associated with a subset of the event space. For example, $H = 0$ is associated with the event TTT , while $H = 1$ is associated with the events $\{HTT, THT, TTH\}$.
- Assuming the probabilities of each of the eight “atomic” events is equal to $\frac{1}{8}$, then the probability mass associated with each value of H is $\{\frac{1}{8}, \frac{3}{8}, \frac{3}{8}, \frac{1}{8}\}$.
- This set of numbers represents the **probability distribution** over H , written $P(H = h) = p_H(h)$. (I will often just write $p(h)$, when the subscript is clear from context.)
- To indicate that the RV (random variable) H is distributed as $p(h)$, we write $H \sim p(h)$.
- The function $p(h)$ is called a probability **mass** function (pmf) if h is discrete, and a probability **density** function (pdf) if h is continuous. In either case, we have $\int_h P(H = h) = 1$ and $\forall h, P(H = h) \geq 0$ for all h in the range of the random variable.
- If we have more than one variable, we can write a joint probability $p_{A,B}(a, b) = P(A = a, B = b)$.
- We can write a **marginal** probability $p_A(a) = \sum_b p_{A,B}(a, b)$.
- Random variables are independent iff $p_{A,B}(a, b) = p_A(a)p_B(b)$.
- We can write a conditional probability as $p_{A|B}(a \mid b) = \frac{p_{A,B}(a,b)}{p_B(b)}$.

Expectations

Sometimes we want the **expectation** of a function, such as $E[g(x)] = \sum_{x \in \mathcal{X}} g(x)p(x)$. Expectations are easiest to think about in terms of probability distributions over discrete events:

- If it is sunny, Marcia will eat three ice creams.
- If it is rainy, she will eat only one ice cream.
- There's a 80% chance it will be sunny.
- The expected number of ice creams she will eat is $0.8 \times 3 + 0.2 \times 1 = 2.6$.

If the random variable X is continuous, the sum becomes an integral:

$$E[g(x)] = \int_{\mathcal{X}} g(x)p(x)dx \quad (1.17)$$

For example, a fast food restaurant in Quebec gives a 1% discount on poutine¹ for every degree below zero. Assuming they use a thermometer with infinite precision, the expected price would be an integral over all possible temperatures,

$$E[\text{price}(x)] = \int_{\mathcal{X}} \min(1, 1 + x) \times \text{original-price} \times p(x)dx. \quad (1.18)$$

(Careful readers will note that the restaurant will apparently pay you for taking poutine, if the temperature falls below -100 degrees celsius.)

1.2 Naïve Bayes

Back to text classification, where we were left wondering how to set the weights θ . Having just reviewed basic probability, we can now take a probabilistic approach to this problem. A Naïve Bayes classifier chooses the weights θ to maximize the *joint* probability of a labeled dataset, $p(\{\mathbf{x}_i, y_i\}_{i \in 1 \dots N})$, where each tuple $\langle \mathbf{x}_i, y_i \rangle$ is a labeled instance.

We first need to define the probability $p(\{\mathbf{x}_i, y_i\}_{i \in 1 \dots N})$. We'll do that through a "generative model," which describes a hypothesized stochastic process that has generated the observed data.²

¹Readers from New Jersey will recognize poutine as a close relative of "disco fries."

²We'll see a lot of different generative models in this course. They are a helpful tool because they clearly and explicitly define the assumptions that underly the form of the probability distribution. For a very readable introduction to generative models in statistics, see Blei (2014).

- For each document i ,
 - draw the label $y_i \sim \text{Categorical}(\mu)$
 - draw the vector of counts $\mathbf{x}_i \sim \text{Multinomial}(\phi_{y_i})$,

The first line of this generative model is “for each document i ”, which tells us to treat each document independently: the probability of the whole dataset is equal to the product of the probabilities of each individual document. The observed word counts and document labels are **independent and identically distributed (IID)**.

$$p(\{\mathbf{x}_i, y_i\}_{i \in 1 \dots N}; \mu, \phi) = \prod_{i=1}^N p(\mathbf{x}_i, y_i; \mu, \phi) \quad (1.19)$$

This means that the words in each document are **conditionally independent** given the parameters μ and ϕ .

The second line indicates $y_i \sim \text{Categorical}(\mu)$, which means that the random variable y_i is a stochastic draw from a categorical distribution with **parameter** μ . A categorical distribution is just like a weighted die: $p_{\text{cat}}(y; \mu) = \mu_y$, where μ_y is the probability of the outcome $Y = y$. For example, if $\mathcal{Y} = \{\text{positive}, \text{negative}, \text{neutral}\}$, we might have $\mu = [0.1, 0.7, 0.2]$. We require $\sum_y \mu_y = 1$ and $\forall_y, \mu_y \geq 0$.

The third and final line invokes the **multinomial distribution**, which is only slightly more complex:

$$p_{\text{mult}}(\mathbf{x}; \phi) = \frac{(\sum_j x_j)!}{\prod_j x_j!} \prod_j \phi_j^{x_j} \quad (1.20)$$

We again require that $\sum_j \phi_j = 1$ and $\forall_j, \phi_j \geq 0$. The second part of the equation is a product over words, with an exponent for each word; recall that $\phi_j^0 = 1$ for all ϕ_j ; this means that the words that have zero count play no role in the overall probability.

The first part of Equation 1.20 doesn't depend on ϕ , and can usually be ignored. Can you see why we need the first part at all?³ We will return to this issue shortly.

We can write $p(\mathbf{x}_i | y_i; \phi)$ to indicate the conditional probability of word counts \mathbf{x}_i given label y_i , with parameter ϕ , which is equal to $p_{\text{mult}}(\mathbf{x}_i; \phi_{y_i})$. By specifying

³Technically, a multinomial distribution requires a second parameter, the total number of counts (the number of words in the document). Even more technically, that number should be treated as a random variable, and drawn from some other distribution. But none of that matters for classification.

the multinomial distribution, we are working with *multinomial naïve Bayes* (MNB). Why “naïve”? Because the multinomial distribution treats each word token independently: the probability mass function factorizes across the counts.⁴ We’ll see this more clearly later, when we show how MNB is an example of linear classification.

Another version of Naïve Bayes

Consider a slight modification to the generative story of NB:

- For each document i
 - Draw the label $y_i \sim \text{Categorical}(\mu)$
 - For each word $n \leq D_i$
 - * Draw the word $w_{i,n} \sim \text{Categorical}(\phi_{y_i})$

This is not quite the same model as multinomial Naive Bayes (MNB): it’s a product of categorical distributions over words, instead of a multinomial distribution over word counts. This means we would generate the words in order, like $p_W(\text{multinomial})p_W(\text{Naive})p_W(\text{Bayes})$. Formally, this is a model for the joint probability $p(\mathbf{w}, y)$, not $p(\mathbf{x}, y)$.

However, as a classifier, it is identical to MNB. The final probabilities are reduced by a factor corresponding to the normalization term in the multinomial, $\frac{(\sum_j x_j)!}{\prod_j x_j!}$. This means that the the probability for a vector of counts \mathbf{x} is larger than the probability for a list of words \mathbf{w} that induces the same counts. But this makes sense: there can be many word sequences that correspond to a single vector counts. For example, *man bites dog* and *dog bites man* correspond to an identical count vector, $\{\text{bites} : 1, \text{dog} : 1, \text{man} : 1\}$, and the total number of word orderings for a given count vector \mathbf{x} is exactly the ratio $\frac{(\sum_j x_j)!}{\prod_j x_j!}$.

From the perspective of classification, none of this matters, because it has nothing to do with the label y or the parameters ϕ . The ratio of probabilities between any two labels y_1 and y_2 will be identical in the two models, as will the maximum likelihood estimates for the parameters μ and ϕ (defined later).

⁴You can plug in any probability distribution to the generative story and it will still be naïve Bayes, as long as you are making the “naïve” assumption that your features are conditionally independent, given the label. For example, a multivariate Gaussian with diagonal covariance would be naïve in exactly the same sense.

Prediction

The Naive Bayes prediction rule is to choose the label y which maximizes $p(\mathbf{x}, y; \phi, \mu)$:

$$\hat{y} = \arg \max_y p(\mathbf{x}, y; \mu, \phi) \quad (1.21)$$

$$= \arg \max_y p(\mathbf{x} \mid y; \phi) p(y; \mu) \quad (1.22)$$

$$= \arg \max_y \log p(\mathbf{x} \mid y; \phi) + \log p(y; \mu) \quad (1.23)$$

Converting to logarithms makes the notation easier. It doesn't change the prediction rule because the log function is monotonically increasing.

Now we can plug in the probability distributions from the generative story.

$$\log p(\mathbf{x}, y; \mu, \phi) = \arg \max_y \log p(\mathbf{x} \mid y; \phi) + \log p(y; \mu) \quad (1.24)$$

$$= \log \left[\frac{(\sum_j x_j)!}{\prod_j x_j!} \prod_j \phi_{y,j}^{x_j} \right] + \log \mu_y \quad (1.25)$$

$$= \log \frac{(\sum_j x_j)!}{\prod_j x_j!} + \sum_j x_j \log \phi_{y,j} + \log \mu_y \quad (1.26)$$

$$\propto \sum_j x_j \log \phi_{y,j} + \log \mu_y \quad (1.27)$$

$$= \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y), \quad (1.28)$$

where

$$\boldsymbol{\theta} = [\boldsymbol{\theta}^{(1)\top}, \boldsymbol{\theta}^{(2)\top}, \dots, \boldsymbol{\theta}^{(K)\top}]^\top \quad (1.29)$$

$$\boldsymbol{\theta}^{(y)} = [\log \phi_{y,1}, \log \phi_{y,2}, \dots, \log \phi_{y,V}, \log \mu_y]^\top \quad (1.30)$$

and $\mathbf{f}(\mathbf{x}, y)$ is a vector of V word counts and an offset, padded by zeros for the labels not equal to y (see equations 1.2-1.5). This ensures that the inner product $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$ only activates the features in $\boldsymbol{\theta}^{(y)}$, which are what we need to compute the joint log-probability $\log p(\mathbf{x}, y)$ for each y . This is a key point: through this notation, we have converted the problem of computing the log-likelihood for a document-label pair $\langle \mathbf{x}_i, y_i \rangle$ into the computation of a vector inner product.

Estimation

The parameters of a multinomial distribution have a simple interpretation: they are the expected frequency for each word. Based on this interpretation, it is tempting to set the parameters empirically, as

$$\phi_{y,j} = \frac{\sum_{i:Y_i=y} x_{i,j}}{\sum_{j'} \sum_{i:Y_i=y} x_{i,j'}} = \frac{\text{count}(y, j)}{\sum_{j'} \text{count}(y, j')} \quad (1.31)$$

This is called a *relative frequency estimator*. It can be justified more rigorously as a *maximum likelihood estimate*.

Our prediction rule in Equation 1.21 is to choose \hat{y} so as to maximize the joint probability $p(\mathbf{x}, y)$. Maximum likelihood estimation proposes to choose the parameters ϕ and μ in much the same way. Specifically, we want to maximize the joint log-likelihood of some **training data**, which consists of a set of annotated examples where we observe both the text and the true label, $\{\mathbf{x}_i, y_i\}_{i \in 1 \dots N}$. Based on the generative model that we have defined, the log-likelihood is:

$$L = \sum_i \log p_{\text{mult}}(\mathbf{x}_i; \phi_{y_i}) + \log p_{\text{cat}}(y_i; \mu). \quad (1.32)$$

Let's continue to focus on the parameters ϕ . Since $p(y)$ is constant in L with respect to these parameters, we can forget it for now,

$$L(\phi) = \sum_i \log p_{\text{mult}}(\mathbf{x}_i; \phi_{y_i}) \quad (1.33)$$

$$= \sum_i \log \frac{(\sum_j x_{i,j})!}{\prod_j x_{i,j}!} \prod_j \phi_{y_i,j}^{x_{i,j}} \quad (1.34)$$

$$= \sum_i \log \left[(\sum_j x_{i,j})! \right] - \sum_j \log (x_{i,j}!) + \sum_j x_{i,j} \log \phi_{y_i,j} \quad (1.35)$$

$$\propto \sum_j x_{i,j} \log \phi_{y_i,j}, \quad (1.36)$$

where I have abused notation by writing \propto to indicate that the left side of Equation 1.36 is equal to the right side plus terms that are constant with respect to ϕ .

We would now like to optimize L , by taking derivatives with respect to ϕ . But before we can do that, we have to deal with a set of constraints:

$$\forall y, \sum_{j=1}^V \phi_{y,j} = 1 \quad (1.37)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

We'll do this by adding a Lagrange multiplier. Solving separately for each label y , we obtain the resulting Lagrangian,

$$\ell[\phi_y] = \sum_{i:Y_i=y} \sum_j x_{ij} \log \phi_{y,j} - \lambda \left(\sum_j \phi_{y,j} - 1 \right) \quad (1.38)$$

We can now differentiate the Lagrangian with respect to the parameter of interest, setting $\frac{\partial \ell}{\partial \phi_{y,j}} = 0$,

$$0 = \sum_{i:Y_i=y} x_{i,j} / \phi_{y,j} - \lambda \quad (1.39)$$

$$\lambda \phi_{y,j} = \sum_{i:Y_i=y} x_{i,j} \quad (1.40)$$

$$\phi_{y,j} \propto \sum_{i:Y_i=y} x_{i,j} = \sum_i \delta(Y_i = y) x_{i,j}, \quad (1.41)$$

where I use two different notations for indicating the same thing: a sum over the word counts for all documents i such that the label $Y_i = y$. This gives a solution for each ϕ_y up to a constant of proportionality. Now recall the constraint $\forall y, \sum_{j=1}^V \phi_{y,j} = 1$; this constraint arises because ϕ_y represents a vector of probabilities for each word in the vocabulary. We can exploit this constraint to obtain an exact solution,

$$\phi_{y,j} = \frac{\sum_{i:Y_i=y} x_{i,j}}{\sum_{j'=1}^V \sum_{i:Y_i=y} x_{i,j'}} \quad (1.42)$$

$$= \frac{\text{count}(y, j)}{\sum_{j'=1}^V \text{count}(y, j')}. \quad (1.43)$$

This is exactly equal to the relative frequency estimator. A similar derivation gives $\mu_y \propto \sum_i \delta(Y_i = y)$, where $\delta(Y_i = y) = 1$ if $Y_i = y$ and 0 otherwise.

Smoothing and MAP estimation

If data is sparse, you may end up with values of $\phi = 0$. For example, the word *Bayesian* may have never appeared in a spam email yet, so the relative frequency estimate $\phi_{\text{SPAM}, \text{Bayesian}} = 0$. But choosing a value of 0 would allow this single feature to completely veto a label, since $P(Y = \text{SPAM} \mid \mathbf{x}) = 0$ if $\mathbf{x}_{\text{Bayesian}} > 0$.

This is undesirable, because it imposes high **variance**: depending on what data happens to be in the training set, we could get vastly different classification rules.

One solution is to **smooth** the probabilities, by adding “pseudo-counts” of α to each count, and then normalizing.

$$\phi_{y,j} = \frac{\alpha + \sum_{i:Y_i=y} x_{i,j}}{\sum_{j'=1}^V \left(\alpha + \sum_{i:Y_i=y} x_{i,j'} \right)} = \frac{\alpha + \text{count}(y, j)}{V\alpha + \sum_{j'=1}^V \text{count}(y, j')} \quad (1.44)$$

This form of smoothing is called “Laplace smoothing”, and it has a nice Bayesian justification, in which we extend the generative story to include ϕ as a random variable (rather than as a parameter). The resulting estimate is called *maximum a posteriori*, or MAP.

Smoothing reduces **variance**, but it takes us away from the maximum likelihood estimate: it imposes a **bias**. In this case, the bias points towards uniform probabilities. Machine learning theory shows that errors on heldout data can be attributed to the sum of bias and variance. Techniques for reducing variance typically increase the bias, so there is a **bias-variance tradeoff**.⁵

- Unbiased classifiers **overfit** the training data, yielding poor performance on unseen data.
- But if we set a very large smoothing value, we can **underfit** instead. In the limit of $\alpha \rightarrow \infty$, we have zero variance: it is the same classifier no matter what data we see! But the bias of such a classifier will be high.
- Navigating this tradeoff is hard. But in general, as you have more data, variance is less of a problem, so you just go for low bias.
- You may wonder if it is possible to choose a separate α_j for each word j , possibly to add larger amounts of smoothing to more common words. Indeed this is possible, and we will talk a great deal about more advanced smoothing techniques in Chapter 5. But I am unaware of any cases where this makes a major positive impact on classification.

Training, testing, and tuning (development) sets

We’ll soon talk about more learning algorithms, but whichever one we apply, we will want to report its accuracy. Really, this is an educated guess about how well the algorithm will do on new data in the future.

To make an estimate of the accuracy, we need to hold out a separate “test set” from the data that we use for estimation (i.e., training, learning). Otherwise, if

⁵The bias-variance tradeoff is covered by Murphy (2012), but see Mohri et al. (2012) for a more formal treatment of this key concept in machine learning theory.

we measure accuracy on the same data that is used for estimation, we will badly overestimate the accuracy that we are likely to get on new data.

Recall that in addition to the parameters μ and ϕ , which are learned on training data, we also have the amount of smoothing, α . This can be considered a “tuning” parameter, and it controls the tradeoff between overfitting and underfitting the training data. Where is the best position on this tradeoff curve? It’s hard to tell in advance. Sometimes it is tempting to see which tuning parameter gives the best performance on the test set, and then report that performance. Resist this temptation! It will also lead to overestimating accuracy on truly unseen future data. For that reason, this is a sure way to get your research paper rejected; in a commercial setting, this mistake may cause you to promise much higher accuracy than you can deliver. Instead, you should split off a piece of your training data, called a “development set” (or “tuning set”).

Sometimes, people average across multiple test sets and/or multiple development sets. One way to do this is to divide your data into “folds,” and allow each fold to be the development set one time. This is called **K-fold cross-validation**. In the extreme, each fold is a single data point. This is called **leave-one-out**.

The Naïvety of Naïve Bayes

Naïve Bayes is simple to work with: estimation and prediction can be done in closed form, and the nice probabilistic interpretation makes it relatively easy to extend the model in various ways. But Naïve Bayes makes assumptions which seriously limit its accuracy, especially in NLP.

- The multinomial distribution assumes that each word is generated independently of all the others (conditioned on the parameter ϕ_y). Formally, we assume conditional independence:

$$p(\text{naïve, Bayes} \mid y) = p(\text{naïve} \mid y)p(\text{Bayes} \mid y). \quad (1.45)$$

- But this is clearly wrong, because words “travel together.” To hone your intuitions about this, try and decide whether you believe

$$p(\text{naïve Bayes}) > p(\text{naïve})p(\text{Bayes}) \quad (1.46)$$

or...

$$p(\text{naïve Bayes}) < p(\text{naïve})p(\text{Bayes}). \quad (1.47)$$

Apply the chain rule!

Traffic lights Dan Klein makes this point with an example about traffic lights. In his hometown of Pittsburgh, there is a $1/7$ chance that the lights will be broken, and both lights will be red. There is a $3/7$ chance that the lights will work, and the north-south lights will be green; there is a $3/7$ chance that the lights work and the east-west lights are green.

The *prior* probability that the lights are broken is $1/7$. If they are broken, the conditional likelihood of each light being red is 1. The prior for them not being broken is $6/7$. If they are not broken, the conditional likelihood of each individual light being red is $1/2$.

Now, suppose you see that both lights are red. According to Naïve Bayes, the probability that the lights are broken is $1/7 \times 1 \times 1 = 1/7 = 4/28$. The probability that the lights are not broken is $6/7 \times 1/2 \times 1/2 = 6/28$. So according to naive Bayes, there is a 60% chance that the lights are not broken!

What went wrong? We have made an independence assumption to factor the probability $P(R, R \mid \text{not-broken}) = P_{\text{north-south}}(R \mid \text{not-broken})P_{\text{east-west}}(R \mid \text{not-broken})$. But this independence assumption is clearly incorrect, because $P(R, R \mid \text{not-broken}) = 0$.

Less Naïve Bayes? Of course we could decide not to make the naive Bayes assumption, and model $P(R, R)$ explicitly. But this idea does not scale when the feature space is large — as it often is in NLP. The number of possible feature configurations grows exponentially, so our ability to estimate accurate parameters will suffer from high variance. With an infinite amount of data, we would be okay; but we never have that. Naïve Bayes accepts some bias, because of the incorrect modeling assumption, in exchange for lower variance.

1.3 Recap

- Documents are represented as “bags of words”, written as the vector \mathbf{x} .
- Feature functions combine the document and the label into a single vector, $\mathbf{f}(\mathbf{x}, y)$.
- Classification can then be performed as a dot-product $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y)$.
- Naive Bayes
 - Define $p(\mathbf{x}, \mathbf{y})$ via a *generative model*
 - Prediction: $\hat{y} = \arg \max_y p(\mathbf{x}_i, y)$

– Learning:

$$\begin{aligned}\theta &= \arg \max_{\theta} p(\mathbf{x}, \mathbf{y}; \theta) \\ p(\mathbf{x}, \mathbf{y}; \theta) &= \prod_i p(x_i, y_i; \theta) = \prod_i p(x_i | y_i) p(y_i) \\ \phi_{y,j} &= \frac{\sum_{i: Y_i=y} x_{ij}}{\sum_{i: Y_i=y} \sum_j x_{ij}} \\ \mu_y &= \frac{\text{count}(Y = y)}{N}\end{aligned}$$

This gives the maximum likelihood estimator (MLE; same as relative frequency estimator)

- The MLE is unbiased, but has high variance. We can navigate the bias-variance tradeoff by adding smoothing pseudo-counts α , reducing variance but adding bias.

Chapter 2

Discriminative learning

2.1 Features

Naïve Bayes is a simple classifier, where the weights are learned based on the joint probability of labels and words. It includes an independence assumption: all features are mutually independent, conditioned on the label.

- We have defined a **feature function** $f(x, y)$, which corresponds to “bag-of-words” features. While these features do violate the independence assumption, the violation is relatively mild.
- We may be interested in other features, which violate independence more severely. Can you think of any?
 - Prefixes, e.g. *anti-*, *im-*, *un-*
 - Punctuation and capitalization
 - Bigrams, e.g. *not good*, *not bad*, *least terrible*, ...

Rich feature sets generally cannot be combined with Naïve Bayes because the distortions resulting from violations of the independence assumption overwhelm the additional power of better features.

$$p(\text{not bad food}|y) \approx p(\text{not}|y)p(\text{bad}|y)p(\text{food}|y) \quad (2.1)$$

$$p(\text{not bad food}|y) \not\approx p(\text{not}|y)p(\text{bad}|y)p(\text{not bad}|y)p(\text{food}|y) \quad (2.2)$$

To use these features, we will need learning algorithms that do not rely on an independence assumption.

2.2 Perceptron

In NB, the weights can be interpreted as parameters of a probabilistic model. But this model requires an independence assumption that usually does not hold, and limits our choice of features. Why not forget about probability and learn the weights in an error-driven way? The perceptron algorithm, shown in Algorithm 1, is one way to do this.

Algorithm 1 Perceptron learning algorithm

```

1: procedure PERCEPTRON( $\mathbf{x}_{1:N}, y_{1:N}$ )
2:   repeat
3:     Select an instance  $i$ 
4:      $\hat{y} \leftarrow \arg \max_y \boldsymbol{\theta}_t^\top \mathbf{f}(\mathbf{x}_i, y)$ 
5:     if  $\hat{y} \neq y_i$  then
6:        $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$ 
7:     else
8:       do nothing
9:   until tired

```

Basically what the algorithm says is this: if you make a mistake, increase the weights for features which are active with the correct label y_i , and decrease the weights for features which are active with the guessed label \hat{y} .

This seems like a cheap heuristic — will it really work? In fact, there is some nice theory for the perceptron.

- If there is a set of weights that correctly separates your data, then your data is **linearly separable**:

$$\forall \mathbf{x}_i, y_i, \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) > \max_{y' \neq y_i} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y'). \quad (2.3)$$

- If your data is linearly separable, it can be proven that the perceptron algorithm will eventually find a separator.
- What if your data is not separable?
 - the number of errors is bounded...
 - but the algorithm will thrash. That is, the weights will cycle between different values, and will never converge.

The perceptron is an **online** learning algorithm.

- This means that it adjusts the weights after every example.
- This is different from Naïve Bayes, which computes corpus statistics and then sets the weights in a single operation. This is a **batch learning** algorithm.
- Other algorithms are **iterative**, in that they perform multiple updates to the weights, but are also **batch**, in that they have to use all the training data to compute the update. We'll mention two of those algorithms later.

Voted (averaged) perceptron

One solution to the thrashing problem is to average the weights across all iterations $t \in 1 \dots T$:

$$\bar{\theta} = \frac{1}{T} \sum_{t=1}^T \theta_t$$

$$\hat{y} = \arg \max_y \bar{\theta}^\top f(x, y)$$

This average will eventually converge, and there is some analysis showing that averaging can improve generalization (Freund and Schapire, 1999; Collins, 2002). However, this rule as described here is not practical. Can you see why not, and how to fix it?

2.3 Loss functions and large-margin classification

Naive Bayes chooses the weights θ by maximizing the joint likelihood $p(\{x_i, y_i\}_i)$. This can be seen, equivalently, as maximizing the log-likelihood (due to the monotonicity of the log function), and as **minimizing** the negative log-likelihood. This negative log-likelihood can therefore be viewed as a **loss function**, which is minimized:

$$\log p(x, y; \theta) = \sum_{i=1}^N \log p(x_i, y_i; \theta) \quad (2.4)$$

$$\ell_{\text{NB}}(\theta; x_i, y_i) = -\log p(x_i, y_i; \theta) \quad (2.5)$$

$$\hat{\theta} = \arg \min_{\theta} \sum_{i=1}^N \ell_{\text{NB}}(\theta, x_i, y_i) \quad (2.6)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

This minimization problem is identical to the maximum-likelihood estimation problem that we solved in the previous chapter. Framing it as minimization may seem confusing and backwards, but loss functions provide a very general framework in which to compare many approaches to machine learning. For example, even though the perceptron is not a probabilistic model, it is also trying to minimize a **loss function**:

$$\ell_{\text{perceptron}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \begin{cases} 0, & y_i = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y) \\ 1, & \text{otherwise} \end{cases} \quad (2.7)$$

The perceptron loss — sometimes called the 0/1 loss — has some pros and cons in comparison with the joint likelihood loss implied by Naive Bayes.

- ℓ_{NB} can suffer **infinite** loss on a single example, which suggests it will overemphasize some examples, and underemphasize others.
- $\ell_{\text{perceptron}}$ treats all errors equally. It only cares if the example is correct, and not about how confident the classifier was. Since we usually evaluate on accuracy, this is a better match.
- $\ell_{\text{perceptron}}$ is non-convex¹ and discontinuous. Finding the global optimum is intractable when the data is not separable.

We can fix this last problem by defining a loss function that behaves more nicely. To do this, let's define the **margin** as

$$\gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \max_{y \neq y_i} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y) \quad (2.8)$$

Then we can write a convex and continuous “hinge loss” as

$$\ell_{\text{hinge}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = \begin{cases} 0, & \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i) \geq 1, \\ 1 - \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i), & \text{otherwise} \end{cases} \quad (2.9)$$

Equivalently, we can write $\ell_{\text{hinge}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = (1 - \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i))_+$, where $(x)_+$ indicates the positive part of x . Essentially, we want a margin of at least 1 between the score for the true label and the best-scoring alternative, which we have written \hat{y} . The hinge and perceptron loss functions are shown in Figure 2.1. Note that the hinge loss upper-bounds the perceptron loss.

¹A function f is convex iff $\alpha f(x_i) + (1 - \alpha)f(x_j) \geq f(\alpha x_i + (1 - \alpha)x_j)$, for all $\alpha \in [0, 1]$ and for all x_i and x_j on the domain of the function. Convexity implies that any local minimum is also a global minimum, and there are effective techniques for optimizing convex functions (Boyd and Vandenberghe, 2004).



Figure 2.1: Hinge and perceptron loss functions

Large-margin online classification

We can write the weight vector $\theta = s\mathbf{u}$, where $\|\mathbf{u}\|_2 = 1$. Think of s as the magnitude and \mathbf{u} as the direction of the vector θ . If the data is separable, there are many values of s which attain zero hinge loss:

$$\gamma(\theta, \mathbf{x}_i, y_i) = \min_{y \neq y_i} \theta^\top \mathbf{f}(\mathbf{x}_i, y_i) - \theta^\top \mathbf{f}(\mathbf{x}_i, y) \quad (2.10)$$

$$= \min_{y \neq y_i} s(\mathbf{u}^\top (\mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, y))) \quad (2.11)$$

If \mathbf{u}^* satisfies $\gamma(s\mathbf{u}^*, \mathbf{x}_i, y_i) > 0$, then there is some smallest value s^* such that $\forall s \geq s^*, \gamma(s\mathbf{u}^*, \mathbf{x}_i, y_i) \geq 1$. This observation suggests that given many possible θ that obtain zero hinge loss, we should choose the one with the smallest norm, since this entails making the least commitment to the training data. This idea underlies the Support Vector Machine (SVM) classifier, which, in its most basic form, solves the optimization problem,

$$\min_{\theta} \|\theta\|_2^2 \quad (2.12)$$

$$s.t. \forall_i \ell_{\text{hinge}}(\theta; \mathbf{x}_i, y_i) = 0. \quad (2.13)$$

In online learning, rather than seeking the feasible θ with the smallest norm, we want to make the smallest magnitude **change** to θ possible. In this way, we hope to limit the thashing problem that we encountered with the perceptron.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Specifically, at each step t , we solve the following optimization problem:

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2 \quad (2.14)$$

$$s.t. \ell_{\text{hinge}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = 0, \quad (2.15)$$

where (\mathbf{x}_i, y_i) is the instance that we draw at step t of the online learning algorithm. This is a constrained quadratic programming problem. Assuming that the constraint can be satisfied (in other words, assuming the problem is linearly separable), the optimal solution is found at,

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \tau_t (\mathbf{f}(y_i, \mathbf{x}_i) - \mathbf{f}(\hat{y}, \mathbf{x}_i)) \quad (2.16)$$

$$\tau_t = \frac{\ell(\boldsymbol{\theta}; \mathbf{x}_i, y_i)}{\|\mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})\|^2}, \quad (2.17)$$

where again \hat{y} is the best scoring y according to $\boldsymbol{\theta}_t$. This solution can be obtained by introducing τ_t as a Lagrange multiplier for the constraint in Equation 2.15.

If the data is not linearly separable, there will be instances for which we can't meet this constraint. To deal with this, we introduce a "slack" variable ξ_i . We use the slack variable to trade off between the constraint (having a large margin) and the objective (having a small change in $\boldsymbol{\theta}$). The tradeoff is controlled by a parameter C .

$$\min w \frac{1}{2} \|\boldsymbol{\theta} - \boldsymbol{\theta}_t\|^2 + C \xi_t \quad (2.18)$$

$$s.t. \ell_{\text{hinge}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) \leq \xi_t, \xi_t \geq 0$$

The solution to Equation 2.18 is,

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \tau_t (\mathbf{f}(y_i, \mathbf{x}_i) - \mathbf{f}(\hat{y}, \mathbf{x}_i)) \quad (2.19)$$

$$\tau_t = \min \left(C, \frac{\ell(\boldsymbol{\theta}; \mathbf{x}_i, y_i)}{\|\mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})\|^2} \right), \quad (2.20)$$

The parameter C is sometimes referred to as the "capacity" of the classifier.

- If C is 0, then infinite slack is permitted, and the weights will never change.
- As $C \rightarrow \infty$, no slack is permitted, and the optimization is identical to Equation 2.15 and Equation 2.17.

(c) Jacob Eisenstein 2014-2015. Work in progress.

This algorithm is called “Passive-Aggressive” (PA; Crammer et al., 2006), because it is passive when the margin constraint is satisfied, but it aggressively changes the weights to satisfy the constraints if necessary.²

- PA is error-driven like the perceptron, but is more stable to violations of separability, like the averaged perceptron.
- PA allows more explicit control than the Averaged Perceptron, due to the C parameter. When C is small, we make very conservative adjustments to θ from each instance, because the slack variables aren’t very expensive. When C is large, we make large adjustments to avoid using the slack variables.
- You can also apply weight averaging to PA.
- **Support vector machines** (SVMs) are another learning algorithm based on the hinge loss (Burges, 1998), but they try to minimize the norm of the weights, rather than the norm of the change in the weights. They are typically trained in **batch** style, meaning that they have to read all the training instances in to compute each update. However, SVMs can also be trained in an online fashion, using a procedure that is quite similar to the Passive-Aggressive online algorithm discussed here (Shalev-Shwartz et al., 2007).

Pros and cons of Perceptron and PA

- Perceptron and PA are error-driven, which means they usually do better in practice than Naïve Bayes.
- They are also online, which means we can learn without having our whole dataset in memory at once. Naïve Bayes can also be estimated online, in the sense that you can stream the data and store the counts.
- The original perceptron doesn’t behave well if the data is not separable, and doesn’t make it easy to control model complexity.
- All these models lack a probabilistic interpretation. Probabilities are useful because they quantify the classification certainty, allowing us to compute expected utility, and to incorporate the classifier in more complex probabilistic models.

²A related algorithm without slack variables is called MIRA, for Margin-Infused Relaxed Algorithm (Crammer and Singer, 2003).

2.4 Logistic regression

Logistic regression is error-driven like the perceptron, but probabilistic like Naïve Bayes. This is useful in case we want to quantify the uncertainty about a classification decision.

Recall that Naïve Bayes selects weights to optimize the joint probability $p(\mathbf{x}, y)$.

- In Naïve Bayes, we factor this as $p(\mathbf{x}, y) = p(\mathbf{x} | y)p(y)$.
- But we could equivalently write $p(\mathbf{x}, y) = p(y | \mathbf{x})p(\mathbf{x})$.

Since we always know \mathbf{x} , we really care only about $p(y | \mathbf{x})$. Logistic regression optimizes this directly. To do this, we have to define the probability function differently. We define the conditional probability directly,

$$p(y | \mathbf{x}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, y'))} \quad (2.21)$$

$$\log p(y | \mathbf{x}) = \sum_{i=1}^N \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y' \in \mathcal{Y}} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y') \quad (2.22)$$

Then the estimation problem is,

$$\hat{\boldsymbol{\theta}} = \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \log p(y_i | \mathbf{x}_i; \boldsymbol{\theta}) \quad (2.23)$$

$$= \arg \max_{\boldsymbol{\theta}} \sum_{i=1}^N \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y' \in \mathcal{Y}} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y') \quad (2.24)$$

Inside the sum, we have the (additive inverse of the) **logistic loss**.

- In binary classification, we can write this as

$$\ell_{\text{logistic}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = -(y_i \boldsymbol{\theta}^\top \mathbf{x}_i - \log(1 + \exp \boldsymbol{\theta}^\top \mathbf{x}_i)) \quad (2.25)$$

- In multi-class classification, we have,³

$$\ell_{\text{logistic}}(\boldsymbol{\theta}; \mathbf{x}_i, y_i) = -(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y' \in \mathcal{Y}} \exp \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}_i, y')) \quad (2.26)$$



Figure 2.2: Hinge, perceptron, and logistic loss functions

The logistic loss is shown in Figure 2.2. Note that logistic loss is also an upper bound on the perceptron loss. Because logistic loss is smooth and convex, we can optimize it through gradient steps,

$$\ell = \sum_{i=1}^N \theta^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y' \in \mathcal{Y}} \exp \theta^\top \mathbf{f}(\mathbf{x}_i, y') \quad (2.27)$$

$$\frac{\partial \ell}{\partial \theta} = \sum_{i=1}^N \mathbf{f}(\mathbf{x}_i, y_i) - \frac{\sum_{y' \in \mathcal{Y}} (\exp \theta^\top \mathbf{f}(\mathbf{x}_i, y')) \mathbf{f}(\mathbf{x}_i, y')}{\sum_{y'' \in \mathcal{Y}} \exp \theta^\top \mathbf{f}(\mathbf{x}_i, y'')} \quad (2.28)$$

$$= \sum_{i=1}^N \mathbf{f}(\mathbf{x}_i, y_i) - \sum_{y' \in \mathcal{Y}} \frac{\exp \theta^\top \mathbf{f}(\mathbf{x}_i, y')}{\sum_{y'' \in \mathcal{Y}} \exp \theta^\top \mathbf{f}(\mathbf{x}_i, y'')} \mathbf{f}(\mathbf{x}_i, y') \quad (2.29)$$

$$= \sum_{i=1}^N \mathbf{f}(\mathbf{x}_i, y_i) - \sum_{y' \in \mathcal{Y}} p(y' | \mathbf{x}_i; \theta) \mathbf{f}(\mathbf{x}_i, y') \quad (2.30)$$

$$= \sum_{i=1}^N \mathbf{f}(\mathbf{x}_i, y_i) - E_{y|x}[\mathbf{f}(\mathbf{x}_i, y)]. \quad (2.31)$$

This gradient has a pleasing interpretation as the difference between the ob-

³The log-sum-exp term is very common in machine learning. It is numerically unstable because you can underflow if the inner product is small, and overflow if the inner product is large. Libraries like `scipy` contain special functions for computing `logsumexp`, but with some thought, you should be able to see how to create an implementation that is numerically stable.

served counts and the expected counts.⁴ Compare this gradient with the perceptron and PA update rules.

The bias-variance tradeoff is handled by penalizing large θ in the objective, adding a term of $\frac{\lambda}{2} \|\theta\|_2^2$. This is called L2 regularization, because of the L2 norm. It can be viewed as placing a zero-mean Gaussian prior distribution on each term of θ , because the log-likelihood under a zero-mean Gaussian is,

$$\log N(\theta_j; 0, \sigma^2) \propto -\frac{1}{2\sigma^2} \theta_j^2, \quad (2.32)$$

so that $\lambda = \frac{1}{\sigma^2}$. This penalty contributes a term of $\lambda\theta$ to the gradient,

$$\ell = \sum_{i=1}^N \theta^\top \mathbf{f}(\mathbf{x}_i, y_i) - \log \sum_{y' \in \mathcal{Y}} \exp \theta^\top \mathbf{f}(\mathbf{x}_i, y') + \frac{\lambda}{2} \|\theta\|_2^2 \quad (2.33)$$

$$\frac{\partial \ell}{\partial \theta} = \sum_{i=1}^N \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y)] - \lambda \theta. \quad (2.34)$$

The effect of this regularizer will cause the estimator to trade off conditional likelihood on the training data for a smaller norm of the weights, and this can help to prevent overfitting. Indeed, regularization is generally considered to be essential to estimating high-dimensional models, as we typically do in NLP. To see why, consider what would happen to the unregularized weight for a base feature j that was active in only one instance \mathbf{x}_i : the conditional likelihood could always be improved by increasing the weight for this feature, so that $\theta_{(j, y_i)} \rightarrow \infty$ and $\theta_{(j, \tilde{y} \neq y_i)} \rightarrow -\infty$, where (j, y) indicates the index of feature associated with $x_{i,j}$ and label y in $\mathbf{f}(\mathbf{x}_i, y)$.

Optimization

In Naive Bayes, the gradient led us to a closed form solution for the parameters θ ; in PA, we obtained a solution for each individual update from a constrained optimization problem. In logistic regression, there are several ways that we could use the gradient of the loss to optimize θ .

Batch optimization In batch optimization, you keep all the data in memory and iterate over it many times. The logistic loss is smooth and convex, so we can find

⁴Recall that the definition of an expected value $E[f(x)] = \sum_x f(x)p(x)$

the global optimum using gradient descent,

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \eta_t \frac{\partial \ell}{\partial \boldsymbol{\theta}}, \quad (2.35)$$

where η_t is some **step size**. In practice, this can be very slow to converge. Second-order (Newton) optimization incorporates the inverse Hessian,

$$H_{i,j} = \frac{\partial^2}{\partial w_i \partial w_j} \ell, \quad (2.36)$$

which leads to much better convergence rates. Unfortunately, in NLP problems, the Hessian matrix (which is quadratic in the number of parameters) is usually too big to deal with.

In practice, people usually apply **quasi-Newton optimization**, which approximates the Hessian matrix, usually L-BFGS (Liu and Nocedal, 1989).⁵ It is usually okay to treat L-BFGS as a black box; you will typically pass it a pointer to a function that computes the likelihood and gradient. In the Python programming language, L-BFGS is provided in `scipy.optimize`.

Online optimization In online optimization, you consider one example (or a “mini-batch” of a few examples) at a time. *Stochastic gradient descent* makes a stochastic online approximation to the overall gradient:

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \eta_t \nabla_{\boldsymbol{\theta}} \ell(\boldsymbol{\theta}^{(t)}, \mathbf{x}, \mathbf{y}) \quad (2.37)$$

$$= \boldsymbol{\theta}^{(t)} - \eta_t (\lambda \boldsymbol{\theta}^{(t)} - \sum_i^N \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y)]) \quad (2.38)$$

$$= (1 - \lambda \eta_t) \boldsymbol{\theta}^{(t)} + \eta_t \sum_i^N \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y)] \quad (2.39)$$

$$\approx (1 - \lambda \eta_t) \boldsymbol{\theta}^{(t)} + N \eta_t (\mathbf{f}(\mathbf{x}_{i(t)}, y_{i(t)}) - E[\mathbf{f}(\mathbf{x}_{i(t)}, y)]) \quad (2.40)$$

where η_t is the **step size** at iteration t , and $\langle \mathbf{x}_{i(t)}, y_{i(t)} \rangle$ is the instance selected at iteration t . (So here we are setting the mini-batch size equal to one.) As always, N is the total number of instances. As above, the expectation is equal to a weighted sum over the labels,

$$E[\mathbf{f}(\mathbf{x}_{i(t)}, y)] = \sum_{y' \in \mathcal{Y}} p(y' | \mathbf{x}_{i(t)}; \boldsymbol{\theta}) \mathbf{f}(\mathbf{x}_{i(t)}, y'). \quad (2.41)$$

⁵A friend told me you can remember the order of the letters as “Large Big Friendly Giants.” Does this help you?

- Note how similar this update is to the perceptron!
- If we set $\eta_t = \eta_0 t^{-\alpha}$ for $\alpha \in [1, 2]$, we have guaranteed convergence.
- We can also just fix η_t to a small value, like 10^{-3} . (This is what we will do in the problem set.)
- In either case, we could tune this parameter on a development set. However, it would be acceptable to just find a value that gives a good regularized log-likelihood on the training set, since this parameter relates to the quality of the optimization, and not the generalization capability of the classifier.
- In theory, we select $\langle \mathbf{x}_{i(t)}, y_{i(t)} \rangle$ at random, but in practice we usually just iterate through the dataset.
- We can fold N into η and λ , so that $\eta^* = N\eta$ and $\lambda^* = \lambda \frac{\eta^*}{N}$. This gives the more compact form,

$$(1 - \lambda^* \eta_t^*) \boldsymbol{\theta}^{(t)} + \eta_t^* (\mathbf{f}(\mathbf{x}_{i(t)}, y_{i(t)}) - E[\mathbf{f}(\mathbf{x}_{i(t)}, y)]) \quad (2.42)$$

For more on stochastic gradient descent, as applied to a number of different learning algorithms, see (Zhang, 2004) and (Bottou, 1998). Murphy (2012) traces SGD to a 1978 paper by GT's own Arkadi Nemirovski (Nemirovski and Yudin, 1978). You can find several recent chapters about online optimization in the edited volume by Sra et al. (2012).

AdaGrad Recent work has shown that you can often learn more quickly by using an **adaptive** step-size, which is different for every feature (Duchi et al., 2011). In the **AdaGrad** algorithm (adaptive gradient), you keep track of the sum of the squares of the gradients for each feature, and rescale the learning rate by its inverse:

$$\mathbf{g}_t = -\mathbf{f}(\mathbf{x}_i, y_i) + \sum_{y' \in \mathcal{Y}} p(y' | \mathbf{x}_i) \mathbf{f}(\mathbf{x}_i, y') + \lambda \boldsymbol{\theta} \quad (2.43)$$

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \frac{\eta}{\sqrt{\sum_{t'=1}^t g_{t,j}^2}} g_{t,j}, \quad (2.44)$$

where j iterates over features in $\mathbf{f}(\mathbf{x}, y)$. The effect of this is that features with consistently large gradients are updated more slowly. Another way to view this update is that rare features are taken more seriously, since their sum of squared gradients will be smaller. AdaGrad seems to require less careful tuning of η , and Dyer (2014) reports that $\eta = 1$ works for a wide range of problems. The

AdaGrad update can apply to any smooth loss function, including the hinge loss defined in Equation 2.9.

Other regularizers

In Equation 2.33, we proposed to **regularize** the estimator of θ by penalizing the squared L2 norm, $\|\theta\|_2^2$. However, this is not the only way to penalize large weights; we might prefer some other norm, such as L_1 or L_0 :

$$L_1 = \|\theta\|_1 = \sum_j |\theta_j| \quad (2.45)$$

$$L_0 = \|\theta\|_0 = \sum_j \delta(\theta_j > 0) \quad (2.46)$$

The L_0 norm penalizes each non-zero weight, so it can be thought of as a form of **feature selection**: optimizing the L_0 -regularized conditional likelihood is equivalent to trading off the log-likelihood against the number of active features. Reducing the number of active features is desirable because the resulting model will be fast, low-memory, and should generalize well, since features that are not very helpful will be pruned away. Unfortunately, the L_0 norm is non-convex and non-differentiable. Optimization under the L_0 norm is NP-hard; indeed, NP-hardness is proven optimization under any L_p norm when $0 \leq p < 1$ (Ge et al., 2011).

However, the L_1 norm is convex, and can be used as an approximation to L_0 (Tibshirani, 1996). Moreover, the L_1 norm also performs feature selection, by driving many of the coefficients to zero; it is therefore known as a **sparsity inducing regularizer**. Gao et al. (2007) compare L_1 and L_2 on a suite of NLP problems, finding that L_1 regularization generally gives similar accuracy to L_2 , but learns models that are between ten and fifty times smaller.

Because the L_1 norm does not have a gradient at $\theta_j = 0$, we must instead optimize using **subgradient** methods. This is only somewhat more complicated in stochastic optimization; Sra et al. (2012) survey approaches for estimation under L_1 and other regularizers.

Other views of logistic regression

Logistic regression is so named because in the binary case where $y \in \{0, 1\}$, we are performing a regression of x against y , after passing the inner product $\theta^\top x$ through a logistic transformation. (You could do a linear regression instead, but this would ignore the fact that the range of y is limited to a few discrete values.)

- Logistic regression is also called **maximum conditional likelihood** (MCL), because it maximizes... the conditional likelihood $p(y | x)$.
- Logistic regression can be viewed as part of a larger family, called **generalized linear models** (GLMs), which include other “link functions” besides the logit, such as the probit function. If you use R, you are probably familiar with `glmnet`, a package for estimating GLMs.
- Logistic regression is also called **maximum entropy**, especially in the earlier NLP literature (Berger et al., 1996). This is due to an alternative formulation, which tries to find the maximum entropy probability function that satisfies moment-matching constraints.

The moment matching constraints specify that the empirical counts of each label-feature pair should match the expected counts:

$$\forall j, \sum_{i=1}^N f_j(\mathbf{x}_i, y_i) = \sum_{i=1}^N \sum_{y \in \mathcal{Y}} p(y | \mathbf{x}_i; \boldsymbol{\theta}) f_j(\mathbf{x}_i, y) \quad (2.47)$$

Note that this constraint will be met exactly when the derivative of the likelihood function (equation 2.31) is equal to zero. However, this will be true for many values of $\boldsymbol{\theta}$. Which should we choose?

The entropy of a conditional likelihood function $p_{Y|X}$ is

$$H(p_{Y|X}) = - \sum_{x \in \mathcal{X}} \tilde{p}_X(x) \sum_{y \in \mathcal{Y}} p_{Y|X}(y | x) \log p_{Y|X}(y | x), \quad (2.48)$$

where $\tilde{p}_X(x)$ is the *empirical probability* of x . We compute an empirical probability by summing over all the instances in training set.

If the entropy is large, this function is smooth across possible values of y ; if it is small, the function is sharp. The entropy is zero if $p(y | x) = 1$ for some particular $Y = y$ and zero for everything else. By saying we want maximum-entropy classifier, we are saying we want to make the weakest commitments possible, while satisfying the moment-matching constraints,

$$\max_{\boldsymbol{\theta}} \quad - \sum_{\mathbf{x}} \tilde{p}_X(\mathbf{x}) \sum_y p_{Y|X}(y | \mathbf{x}; \boldsymbol{\theta}) \log p_{Y|X}(y | \mathbf{x}; \boldsymbol{\theta}) \quad (2.49)$$

$$s.t. \quad \forall j, \sum_{i=1}^N f_j(\mathbf{x}_i, y_i) = \sum_{i=1}^N \sum_y p_{Y|X}(y | \mathbf{x}_i; \boldsymbol{\theta}) f_j(\mathbf{x}_i, y). \quad (2.50)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

The solution to this constrained optimization problem is identical to the maximum conditional likelihood (logistic-loss) formulation we considered in the previous section.

This view of logistic regression is arguably a little dated, but it is useful to understand, especially when reading classic papers from the 1990s. The information-theoretic concept of entropy will pop up again a few times in the course. For a tutorial on maximum entropy, see <http://www.cs.cmu.edu/afs/cs/user/abberger/www/html/tutorial/tutorial.html>.

2.5 Summary of learning algorithms

- **Naive Bayes.** pros: easy and probabilistic. cons: arguably optimizes wrong objective; usually has poor accuracy, especially with overlapping features.
- **Perceptron and PA.** pros: easy, online, and error-driven. cons: not probabilistic. this can be bad in pipeline architectures, where the output of one system becomes the input for another.
- **Logistic regression.** pros: error-driven and probabilistic. cons: batch learning requires black-box software; hinge loss sometimes yields better accuracy than logistic loss.

For more details, see Table 2.1.

What about non-linear classification?

The feature spaces that we consider in NLP are usually huge, so non-linear classification can be quite difficult. When the feature dimension V is larger than the number of instances N — often the case in NLP — you can always learn a linear classifier that will perfectly classify your training instances.⁶ This makes selecting an appropriate **non-linear** classifier especially difficult. Nonetheless, there are some approaches to non-linear learning in NLP:

- You can add **features**, such as bigrams, which are non-linear combinations of other features. For example, the base feature $\langle \text{coffee house} \rangle$ will not fire unless both features $\langle \text{coffee} \rangle$ and $\langle \text{house} \rangle$ also fire.
- Another option is to apply non-linear transformations to the feature vector. Recall that the feature function $f(x, y)$ may be composed of a vector of

⁶Assuming your feature matrix is full-rank.

word counts, padded by zeros. We can think of these word counts as basic features, and apply non-linear transformations, such as $x \circ x$ or $|x|$.

- There is some work in NLP on using kernels for strings, bags-of-words, sequences, trees, etc (Collins and Duffy, 2001; Zelenko et al., 2003). Kernel-based learning can be seen as a generalization of algorithms such k -nearest-neighbors, which classifies instances by considering the labels of the k most similar instances in the training set (Hastie et al., 2009). This is beyond the scope of this text, but see Murphy (2012) for more details.
- Boosting (Freund et al., 1999) and decision tree algorithms (Schmid, 1994) sometimes do well on NLP tasks, but they are used less frequently these days, especially as the field increasingly emphasizes big data and simple classifiers.
- More recent work has shown how **deep learning** can perform non-linear classification. One way to use deep learning in NLP is by learning word representations, while jointly learning how these representations combine to classify instances (Collobert and Weston, 2008). This approach is very hot at the moment, so I will discuss it towards the end of the semester.

	Naive Bayes	Logistic Regression	Perceptron	PA
Objective	Joint likelihood	Conditional likelihood	0/1 loss	Hinge loss
estimation	$\max \sum_i \log \mathbf{p}(\mathbf{x}_i, y_i)$	$\max \sum_i \log \mathbf{p}(y_i \mathbf{x}_i)$	$\min \sum_i \delta(y_i, \hat{y})$	$\sum_i [1 - \gamma(\boldsymbol{\theta}; \mathbf{x}_i, y_i)] +$
tuning	$\theta_{ij} = \frac{c(\mathbf{x}_i, y=j) + \alpha}{c(y=j) + V\alpha}$	$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \sum_i \mathbf{f}(\mathbf{x}_i, y_i) - E[\mathbf{f}(\mathbf{x}_i, y)]$	$\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y})$	$\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \tau_k(\mathbf{f}(\mathbf{x}_i, y_i) - \mathbf{f}(\mathbf{x}_i, \hat{y}))$
complexity	smoothing α	regularizer $\lambda \ \boldsymbol{\theta}\ _2^2$	weight averaging	slack penalty C
easy?	$\mathcal{O}(NV)$	$\mathcal{O}(NVT)$	$\mathcal{O}(NVT)$	$\mathcal{O}(NVT)$
probabilities?	very	not really	yes	yes
	yes	yes	no	no
features?	no	yes	yes	yes

Table 2.1: Comparison of classifiers. N = number of examples, V = number of features, T = number of instances.

(leave this page blank or the next page gets messed up)

Chapter 3

Linguistic applications of classification

Having learned some techniques for classification, let's now see how they can be applied to typical problems in natural language technology.

3.1 Sentiment and opinion analysis

A popular NLP technology is automatically determining the “sentiment” or “opinion polarity” of documents such as product reviews and social media posts. For example, marketers are interested to know how people respond to advertisements, services, and products (Hu and Liu, 2004); social scientists are interested in how peoples’ emotions are affected by phenomena such as the weather (Hannak et al., 2012), and how emotions spread over social networks (Coviello et al., 2014; Miller et al., 2011). In the field of **digital humanities**, literary scholars track plot structures through the flow of sentiment across a novel (Jockers, 2015). A comprehensive analysis of this broad literature is beyond the scope of this chapter, but see survey manuscripts by Pang and Lee (2008) and Liu (2012, 2015).

Sentiment analysis can be framed as a fairly direct application of document classification, as long as reliable labels can be obtained. In the simplest case, sentiment analysis can be treated as a two or three-class problem, with sentiments of POSITIVE, NEGATIVE, and possibly NEUTRAL. Such annotations could be annotated by hand, or obtained automatically through a variety of means:

- Tweets containing happy emoticons can be marked as positive, sad emoticons as negative (Read, 2005; Pak and Paroubek, 2010).

- Reviews with four or more stars can be marked as positive, two or fewer stars as negative (Pang et al., 2002).
- Statements from politicians who are voting **for** a given bill are marked as positive (towards that bill); statements from politicians voting against the bill are marked as negative (Thomas et al., 2006).

After obtaining the annotations, several design decisions may be taken in construction of the feature vector $f(x, y)$:

Preprocessing One question is whether the vocabulary should be case sensitive: do we distinguish *great*, *Great*, and *GREAT*? What about *cooooooooool*? In social media text, this sort of **expressive lengthening** can cause the vocabulary size to explode (Brody and Diakopoulos, 2011); we might want to somehow **normalize** the text (Sproat et al., 2001) to collapse the vocabulary again.

A related issue is that suffixes may be irrelevant to the sentiment orientation of each word: for example, *love*, *loved*, and *loving* are all positive, so perhaps we should eliminate the suffix and group them together. The removal of these suffixes is called **stemming** when it is done at the character level (leaving roots like *lov-*), and is called **lemmatization** when the goal is to identify the underlying base word (in this case, *love*). Both of these methods will be discussed in detail in chapter 6 and chapter 7.

Still another preprocessing decision involves **tokenization**: breaking the text into tokens. This is more complicated than simply looking for whitespace, since we may want to tokenize items such as *well-bred* into $\langle \text{well}, \text{bred} \rangle$, *isn't* into $\langle \text{is}, \text{n't} \rangle$; at the same time, we would like to keep *U.S.* as a single token. This too will be discussed in chapter 7.

Vocabulary In some cases, it is preferable not to include all words in the vocabulary. Words such as *the*, *to*, and *and* seem intuitively to play little role in expressing sentiment or opinion, yet they are very frequent; removing these **stopwords** may therefore improve the classifier. This is typically done by creating a list and simply matching all items on the list. More aggressively, we might assume that sentiment is typically carried by **adjectives** and **adverbs** (see Chapter 8), and therefore we could focus on these words (Hatzivassiloglou and McKeown, 1997; Turney, 2003). However, Pang et al. (2002) find that in their case, eliminating non-adjectives causes the performance of the classifier to decrease.

Count or binary? Finally, we may consider whether we want our feature vector to include the **count** of each word, or its mere **presence**. This gets at a subtle limitation of linear classification: two *failures* may be worse than one, but is it really twice as bad? A more flexible classifier could assign diminishing weight to each additional instance, but this is hard to do in the linear classification framework, and it's hard to see how much the weight should diminish. Pang et al. (2002) take a simpler approach, using binary presence/absence indicators in the feature vector: $f_i(\mathbf{x}, y) \in \{0, 1\}, \forall i$. They find that classifiers trained on these binary feature vectors outperform classifiers trained on count-based features.

A more challenging version of opinion analysis is to determine not just the class of a review, but its rating on a numerical scale (Pang and Lee, 2005). If the scale is continuous, we might take a regression approach, identifying a set of weights θ so as to minimize the squared error of a predictor $\hat{y} = \theta^\top \mathbf{x} + b$, where b is an offset. We can remove the offset by adding a feature to \mathbf{x} whose value is always 1; the corresponding weight in θ is then equivalent to b . Least squares regularization has a closed form solution,

$$\theta = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}, \quad (3.1)$$

where \mathbf{y} is a column vector of size N , containing all ratings in the training data, and \mathbf{X} is an $N \times D$ matrix containing all D features for all N instances. If we place an L2 regularizer on θ , with penalty $\lambda \|\theta\|_2^2$, the resulting problem is called **ridge regression**. It too has a closed form solution,

$$\theta = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbb{I})^{-1} \mathbf{X}^\top \mathbf{y}. \quad (3.2)$$

If the rating scale is discrete, $y \in \{1, 2, \dots, K\}$, we can take a **ranking** approach (Crammer and Singer, 2001), in which scores $\theta^\top \mathbf{x}$ are discretized into ranks, by also learning a set of boundaries, $b_0 = -\infty \leq b_1 \leq \dots \leq b_K$. The learning algorithm consists in making perceptron-like updates to both θ and \mathbf{b} . This approach is ideal for settings like predicting a 1-10 rating or a grade (A - F); instead of learning one vector θ for every rank, we can learn a single θ , and then just partition the output space.

[**todo: Other topics to cover:**]

- subjectivity
- sentence-level versus document-level sentiment
- negation and the role of syntax

(c) Jacob Eisenstein 2014-2015. Work in progress.

- targeted sentiment
- Stance classification

3.2 Word sense disambiguation

Consider the the following headlines:

- (3.1) *Iraqi head seeks arms*
- (3.2) *Prostitutes appeal to Pope*
- (3.3) *Drunk gets nine years in violin case*¹

They are ambiguous because they contain words that have multiple meanings, or **senses**. Word Sense Disambiguation (WSD) is the problem of identifying the intended sense of each word token in a document. WSD is part of a larger field of research called **lexical semantics**, which is concerned with meanings of the words.

Problem definition

Part-of-speech ambiguity (e.g., noun versus verb, as in *she is **heading** out of town*) is usually considered to be a different problem from WSD. Here we are focusing on ambiguity between senses that are all the same part-of-speech, and in part-of-speech tagging evaluations, it is often assumed that the correct part-of-speech has already been identified. [todo: why?] From a linguistic perspective, senses are not really properties of words, but of **lemmas**, which are groups of inflected forms, e.g. (*arm/N, arms/N*), (*arm/V, arms/V, armed/V, arming/V*), where *arm/N* indicates the word *arm* tagged as a noun (*V* is for verb). So the WSD problem can be defined as identifying the correct sense for each word token from an inventory associated for the word's lemma.

How many word senses?

Words (lemmas) may have *many* more than two senses. For example, the word *serve* would seem to have at least the following senses:

- [FUNCTION]: *The tree stump served as a table*

¹These examples, and many more, can be found at <http://www.ling.upenn.edu/~beatrice/humor/headlines.html>

- [ENABLE]: *His evasive replies only served to heighten suspicion*
- [DISH]: *We serve only the rawest fish here*
- [ENLIST]: *She served her country in the marines*
- [JAIL]: *He served six years in Alcatraz*
- [TENNIS]: *Nobody can return his double-reverse spin serve*
- [LEGAL]: *They were served with subpoenas²*

How do we know that these senses are really different? Linguists often design tests for this purpose, and one such test is to construct a **zeugma**, which combines antagonistic senses in an uncomfortable way:

(3.4) *Which flight serves breakfast?*

(3.5) *Which flights serve Tuscon?*

(3.6) **Which flights serve breakfast and Tuscon?³*

The asterisk is a linguistic notation for utterances which would not be judged to be grammatical by fluent speakers of a language. To the extent that you think that (3.6) is ungrammatical, you should agree that (3.4) and (3.5) refer to distinct senses of the lemma *serve*.

The WSD task: Output What should the output of WSD be? What are the possible senses for each word? We could just look in the dictionary. But rather than using a traditional dictionary, WSD research is dominated by a computational resource called WORDNET (<http://wordnet.princeton.edu>). WordNet is organized in terms of lemmas rather than words. An example of a wordnet entry is shown in Figure 3.1

WordNet consists of roughly 100,000 **synsets**, groups of words or phrases with an identical meaning. (e.g., {CHUMP¹, FOOL², SUCKER¹, MARK⁹}). A lemma is **polysemous** if it participates in multiple synsets. Besides **synonymy**, WordNet also describes many other lexical relationships, including:

antonymy *x* means the opposite of *y*, e.g. FRIEND-ENEMY;

hyponymy *x* is a special case of *y*, e.g. RED-COLOR; the inverse relationship is **hypernymy**;

²Examples from Dan Klein's lecture notes, [http://www.cs.berkeley.edu/~klein/cs294-7/SP07%20cs294%20lecture%205%20--%20maximum%20entropy%20\(6pp\).pdf](http://www.cs.berkeley.edu/~klein/cs294-7/SP07%20cs294%20lecture%205%20--%20maximum%20entropy%20(6pp).pdf)

³I believe this example is from Jurafsky and Martin (2009) [**todo: but check**].



Figure 3.1: Example wordnet entry, from <http://wordnet.princeton.edu>

meronymy x is a part of y , e.g., WHEEL-BICYCLE; the inverse relationship is **holonymy**.

WordNet has played a big role in helping WSD move from toy systems to large-scale quantitative evaluations. However, some have argued that WordNet's sense granularity is too fine (Ide and Wilks, 2006); more fundamentally, the premise that word senses can be differentiated in a task-neutral way has been criticized as linguistically naïve (Kilgarriff, 1997). One way of testing this question is to ask whether people tend to agree on the appropriate sense for example sentences: according to Mihalcea et al. (2004), humans agree on roughly 70% of examples using WordNet senses; far better than chance, but perhaps less than we might like.

A range of tasks have been proposed for WSD:

- **Synthetic** data: different words are conflated (*banana-phone*), the system must identify the original word.
- **Lexical sample**: disambiguate a few target words (e.g., *plant* etc). This is what was used in the first large-scale WSD evaluation, SENSEVAL-1 (1998).[todo: citation]
- **All-words** WSD: a sense must be identified for every token.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- A **semantic concordance** is a corpus in which each open-class word (nouns, verbs, adjectives, and adverbs) is tagged with its word sense from the target dictionary or thesaurus.
- SEMCOR is a semantic concordance built from 234K tokens of the Brown corpus.

As of Sunday_n¹ night_n¹ there was_v⁴ no word_n² ...

WSD as Classification

So, how can we tell living *plants* from manufacturing *plants*? The key information often lies in the **context**:

- (3.7) *Town officials are hoping to attract new manufacturing plants through weakened environmental regulations.*
- (3.8) *The endangered plant plays an important role in the local ecosystem.*

Bag-of-words models are a very typical approach. For example,

$$f(y, \text{bank}, I \text{ went to the bank to deposit my paycheck}) = \{ \langle \text{went}, y \rangle : 1, \langle \text{deposit}, y \rangle : 1, \langle \text{paycheck}, y \rangle : 1 \}$$

Some examples:⁴

- *bank*[FINANCIAL]:

a an and are ATM Bonnie card charges check Clyde criminals deposit famous for get I much My new overdraft really robbers the they think to too two went were

- *bank*[RIVER]:

a an and big campus cant catfish East got grandfather great has his I in is Minnesota Mississippi muddy My of on planted pole pretty right River The the there University walk Wets

An extension of bag-of-words models is to encode the position of each context word, e.g.,

$$f(y, \text{bank}, I \text{ went to the bank to deposit my paycheck}) = \{ \langle i - 3, \text{went}, y \rangle : 1, \langle i + 2, \text{deposit}, y \rangle : 1, \langle i + 4, \text{paycheck}, y \rangle : 1 \}$$

⁴todo: reconcile with examples above

Jurafsky and Martin (2009) call these **collocation features**. Other approaches include more information about the sentence structure, such as the part-of-speech tag for each word, and the words with which it is syntactically linked in the sentence (see chapter 12).

After deciding on the features, we can train a classifier to predict the right sense of each word — assuming enough labeled examples can be accumulated. This is difficult, because each polysemous lemma requires its own training set: having a good classifier for *bank* is of no help at all towards disambiguating *plant*. For this reason, **unsupervised** and **semisupervised** methods are particularly popular for WSD (Yarowsky, 1995). We will talk about related methods in chapter 4 and chapter 20. Unsupervised methods typically lean heavily on the heuristic “one sense per discourse”, meaning roughly that a lemma will have a consistent sense throughout any given document. Based on this heuristic, we can propagate information from high-confidence instances to lower-confidence instances in the same document. For a survey on word sense disambiguation, see Navigli (2009).

3.3 Other applications

- Author identification
- Author demographics, maybe
- Language classification

Chapter 4

Learning without supervision

So far we've assumed the following setup:

- A **training set** where you get observations x_i and labels y_i
- A **test set** where you only get observations x_i

What if you never get labels y_i ? For example, you get a bunch of text, and you suspect that there are at least two different meanings for the word *concern*.¹

As described in chapter 3, in supervised word sense disambiguation, we often build feature vectors from the words that appear in the context of the word that we are trying to disambiguate. For example, for the word *concern*, the immediate context might typically include words from one of the following two groups:

1. *services, produces, banking, pharmaceutical, energy, electronics*
2. *about, said, that, over, in, with, had*

Now suppose we were to scatterplot each instance of *concern* on a graph, so that the x-axis is the density of words in group 1, and the y-axis is the density of words in group 2. In such a graph, shown in Figure 4.1, two or more blobs might emerge. These blobs would correspond to the different sense of *concern*.

But in reality, we don't know the word groupings in advance.² We have to try to apply the same idea in a very high dimensional space, where every word gets its own dimension — and most dimensions are irrelevant!

¹The example is from Pedersen and Bruce (1997).

²One approach, which we do not consider here, would be to get them from some existing resource, such as the dictionary definition. (Lesk, 1986)

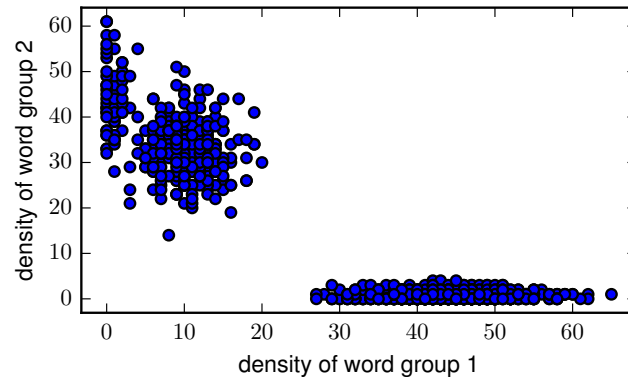


Figure 4.1: Counts of words from two different context groups

Now here’s a related scenario, from a different problem. Suppose you download thousands of news articles, and make a scatterplot, where each point corresponds to a document: the x-axis is the frequency of the word *hurricane*, and the y-axis is the frequency of the word *election*. Again, three clumps might emerge: one for documents that are largely about the hurricane, another for documents largely about the election, and a third clump for documents about neither topic.

These examples are intended to show that we can find structure in data, even without labels — just look for clumps in the scatterplot of features. But again, in reality we cannot make scatterplots of just two words; we may have to consider hundreds or thousands of words. It would be impossible to visualize such a high-dimensional scatterplot, so we will need to design algorithmic approaches to finding the structure (in this case, the “blobs”).

4.1 K-means clustering

You might know about classic clustering algorithms like *K*-means. These algorithms are iterative:

1. Guess the location of cluster centers.
2. Until converged:
 - a) Assign each point to the nearest center.
 - b) Re-estimate the centers as the mean of the assigned points.

(c) Jacob Eisenstein 2014-2015. Work in progress.

This can be viewed as an algorithm for finding coherent “blobs” of documents in high-dimensional data. When we assign each point to its nearest center, we are choosing which blob it is in; when we re-estimate the location of the centers, we are determining the defining characteristic of each blob. K -means is a classic algorithmic that has been used and modified in thousands of papers (Jain, 2010); for an application of K -means to word sense induction, see Pantel and Lin (2002).

Of the many variants of K -means, one that is particularly relevant for our purposes is called “soft” K -means. The key difference is that instead of directly assigning each point \mathbf{x}_i to a specific cluster z_i , we assign it a **distribution** over clusters $q_i(z_i)$, so that $\sum_k q_i(k) = 1$, and $\forall_k 0 \leq q_i(k) \leq 1$. The centroid of each cluster is then computed from a **weighted average** of the points in the cluster, where the weights are taken from the q distribution.

We will now explore a more principled, statistical version of soft K -means, called EM clustering. By understanding the statistical principles underlying the algorithm, we can extend it in a number of ways.

4.2 The Expectation-Maximization (EM) Algorithm

Let’s go back to the Naive Bayes model:

$$\log p(\mathbf{x}, \mathbf{y}; \phi, \mu) = \sum_i \log p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \quad (4.1)$$

For example, \mathbf{x} can describe the documents that we see today, and \mathbf{y} can correspond to their labels. But suppose we never observe y_i ? Can we still do anything with this model?

Since we don’t know \mathbf{y} , let’s marginalize it:

$$\log p(\mathbf{x}) = \log \sum_{\mathbf{y}} p(\mathbf{x} | \mathbf{y}; \phi) p(\mathbf{y}; \mu) \quad (4.2)$$

$$= \log \sum_{\mathbf{y}} \prod_i p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \quad (4.3)$$

$$= \sum_i \log \sum_{y_i} p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \quad (4.4)$$

We are now concerned with maximizing the likelihood of $\mathbf{x}_{1:N}$, which is our (unlabeled) observed data. Why is this a good thing to maximize? If we don’t have

labels, discriminative learning is impossible (there's nothing to discriminate), so maximum likelihood is all we have.³

To estimate this model, we introduce an **auxiliary variable** q_i , for each y_i . We want q_i to be a **distribution**, so we have the usual constraints: $\sum_y q_i(y) = 1$ and $\forall y, q_i(y) \geq 0$. In other words, q_i defines a probability distribution over \mathcal{Y} , for each instance i .

Now since $\frac{q_i(y)}{q_i(y)} = 1$,

$$\log p(\mathbf{x}) = \sum_i \log \sum_{y_i} p(\mathbf{x}_i | y_i; \phi) p(y_i; \mu) \frac{q_i(y)}{q_i(y)} \quad (4.5)$$

$$= \sum_i \log E_q \left[\frac{p(\mathbf{x}_i | y; \phi) p(y; \mu)}{q_i(y)} \right], \quad (4.6)$$

by the definition of expectation, $E_q f(x) = \sum_q q(x) f(x)$. Note that E_q just means the expectation under the distribution q .

Now we apply **Jensen's inequality**, which says that because \log is concave, we can push it inside the expectation, and obtain a lower bound.

$$\log p(\mathbf{x}) \geq \sum_i E_q \left[\log \frac{p(\mathbf{x}_i | y; \phi) p(y; \mu)}{q_i(y)} \right] \quad (4.7)$$

$$\mathcal{J} = \sum_i E_q [\log p(\mathbf{x}_i | y; \phi)] + E_q [\log p(y; \mu)] - E_q [\log q_i(y)] \quad (4.8)$$

By maximizing \mathcal{J} , we are maximizing a lower bound on the joint log-likelihood $\log p(\mathbf{x})$. Now, \mathcal{J} is a function of two arguments:

- the distributions $q_i(\mathbf{y})$ for each i
- the parameters μ and ϕ

We'll optimize with respect to each of these in turn, holding the other one fixed.

³This formulation of $p(\mathbf{x})$ is sometimes called a **mixture model**, since each instance is modeled as a latent mixture over components $\phi_1 \dots \phi_K$.

The E-step

First, we expand the expectation in the lower bound as:

$$\mathcal{J} = \sum_i E_q[\log p(\mathbf{x}_i | y; \phi)] + E_q[\log p(y; \mu)] - E_q[\log q_i(y)] \quad (4.9)$$

$$= \sum_i \sum_y q_i(y) (\log p(\mathbf{x}_i | y; \phi) + \log p(y; \mu) - \log q_i(y)) \quad (4.10)$$

As in Naïve Bayes, we have a “sum-to-one” constraint: in this case, $\sum_y q_i(y) = 1$. Once again, we incorporate this constraint into a Lagrangian:

$$\mathcal{J}_q = \sum_i^N \sum_{y \in \mathcal{Y}} q_i(y) (\log p(\mathbf{x}_i | y; \phi) + \log p(y; \mu) - \log q_i(y)) + \lambda_i (1 - \sum_y q_i(y)) \quad (4.11)$$

$$\frac{\partial \mathcal{J}_q}{\partial q_i(y)} = \log p(\mathbf{x}_i | y; \phi) + \log p(y; \mu) - \log q_i(y) - 1 - \lambda_i \quad (4.12)$$

$$\log q_i(y) = \log p(\mathbf{x}_i | y; \phi) + \log p(y; \mu) - 1 - \lambda_i \quad (4.13)$$

$$q_i(y) \propto p(\mathbf{x}_i | y; \phi) p(y; \mu) \quad (4.14)$$

$$\propto p(\mathbf{x}_i, y; \phi, \mu) \quad (4.15)$$

Since q_i is defined over the labels \mathcal{Y} , we normalize it as,

$$q_i(y) = \frac{p(\mathbf{x}_i, y; \phi, \mu)}{\sum_{y' \in \mathcal{Y}} p(\mathbf{x}_i, y'; \phi, \mu)} \quad (4.16)$$

$$= p(y | \mathbf{x}_i; \phi, \mu) \quad (4.17)$$

After normalizing, each $q_i(y)$ — which is the soft distribution over clusters for data \mathbf{x}_i — is set to the posterior probability $p(y | \mathbf{x}_i)$ under the current parameters μ, ϕ . This is called the E-step, or “expectation step,” because it is derived from updating the bound on the expected likelihood under $q(y)$.

The M-step

Next, we hold $q(y)$ fixed and maximize the bound with respect to the parameters, ϕ and μ . Lets focus on ϕ , which parametrizes the likelihood, $p(\mathbf{x} | y; \phi)$. Again,

(c) Jacob Eisenstein 2014-2015. Work in progress.

we have a constraint that $\sum_j \phi_{y,j} = 1$, so we start by forming a Lagrangian,

$$\mathcal{J}_\phi = \sum_i^N \sum_{y \in \mathcal{Y}} q_i(y) (\log p(\mathbf{x}_i | y; \phi) + \log p(y; \mu) - \log q_i(y)) + \sum_{y \in \mathcal{Y}} \lambda_y (1 - \sum_j^V \phi_{y,j}) \quad (4.18)$$

$$\frac{\partial \mathcal{J}_\phi}{\partial \phi_{y,j}} = \sum_i^N q_i(y) \frac{x_{i,j}}{\phi_{y,j}} - \lambda_y \quad (4.19)$$

$$\lambda_y \phi_{y,j} = \sum_i^N q_i(y) x_{i,j} \quad (4.20)$$

$$\phi_{y,j} \propto \sum_i^N q_i(y) x_{i,j} \quad (4.21)$$

Now because $\sum_j \phi_{y,j} = 1$, we can normalize as follows,

$$\phi_{y,j} = \frac{\sum_i^N q_i(y) x_{i,j}}{\sum_{j' < V} \sum_i^N q_i(y) x_{i,j'}} \quad (4.22)$$

$$= \frac{E_q[\text{count}(y, j)]}{E_q[\text{count}(y)]}. \quad (4.23)$$

So ϕ_y is now equal to the relative frequency estimate of the **expected counts** under the distribution $q(y)$.

- As in supervised Naïve Bayes, we can apply smoothing to add α to all these counts.
- The update for μ is identical: $\mu_y \propto \sum_i q_i(y)$, the expected proportion of cluster $Y = y$. If needed, we can add smoothing here too.
- So, everything in the M-step is just like Naive Bayes, except we used expected counts rather than observed counts.

In some cases, there is no closed form solution for the parameters in the M-step. We may therefore run gradient-based optimization at each M-step, or we may simply take a single step along the gradient step and then return to the E-step (Berg-Kirkpatrick et al., 2010).

(c) Jacob Eisenstein 2014-2015. Work in progress.

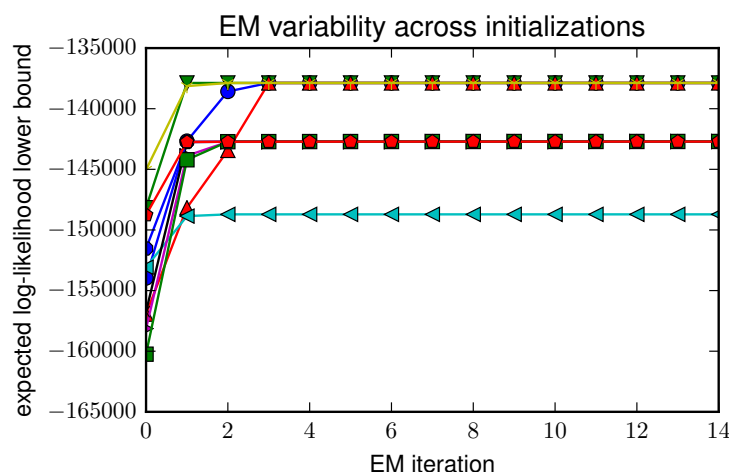


Figure 4.2: Sensitivity of expectation maximization to initialization

Coordinate ascent

Algorithms that alternate between updating various subsets of the parameters are called “coordinate ascent” algorithms.

The objective function \mathcal{J} is **biconvex**, meaning that it is separately convex in $q(\mathbf{y})$ and $\langle \mu, \phi \rangle$, but it is not jointly convex in all three terms. In the coordinate ascent algorithm that we have defined, each step is guaranteed not to decrease \mathcal{J} . This is sometimes called “hill climbing”, because you never go down. Specifically, EM is guaranteed to converge to a **local optima** — a point which is as good or better than any of its immediate neighbors. But there may be many such points, and the overall procedure is **not** guaranteed to find a global maximum. Figure 4.2 shows the objective function for EM with ten different random initializations: while the objective function increases monotonically in each run, it converges to several different values.

The fact that there is no guarantee of global optimality means that initialization is important: where you start can determine where you finish. This is not true in most of the supervised learning algorithms that we have considered, such as logistic regression; in that case, we are optimizing $\log p(\mathbf{y} \mid x; \theta)$, which by construction is convex with respect to the parameter θ . So, for logistic regression — and for many other supervised learning algorithms — we don’t need to worry about initialization, because it won’t affect our ultimate solution: we are guaranteed to reach the global minimum. Recent work on **spectral learning** has sought

to obtain similar guarantees for “latent variable” models, such as the case we are considering now, where x is observed and y is latent. This work is briefly touched on in section 4.4.

Variants In “hard EM”, each $q_i(y)$ distribution assigns probability of 1 to a single \hat{y}_i , and probability of 0 to all others (Neal and Hinton, 1998). This is similar in spirit to K-means clustering; indeed, if the likelihood $p(x | y)$ is Gaussian, then hard EM is identical to K-means clustering with a Euclidean distance function. In problems where the space \mathcal{Y} is large, it can be easier to find the maximum likelihood value \hat{y} than it is to compute the entire distribution $q_i(y)$. Spitkovsky et al. (2010) show that this hard EM can outperform standard EM in some cases. Hard EM can be generalized by adding an additional term that penalizes the **entropy** of q , $H(q_i) = -\sum_y q_i(y) \log q_i(y)$ (see chapter 5 for much more on entropy), yielding a range of variants of the EM algorithm (Samdani et al., 2012).

Another variant of the coordinate ascent procedure combines EM with stochastic gradient descent (SGD). In this case, we can do a local E-step at each instance i , and then immediately make a gradient update to the parameters $\langle \mu, \phi \rangle$. This is particularly relevant in cases where there is no closed-form solution for the parameters, so that gradient ascent will be necessary in any case. This algorithm is called “incremental EM” by Neal and Hinton (1998), and online EM by Sato and Ishii (2000) and Cappé and Moulines (2009). Liang and Klein (2009) apply a range of different online EM variants to NLP problems, obtaining better results than standard EM in many cases.

How many clusters?

All along, we have assumed that the number of clusters $K = \#|\mathcal{Y}|$ is given. In some cases, this assumption is valid. For example, the dictionary or WordNet might tell us the number of senses for a word. In other cases, the number of clusters should be a tunable parameter: some readers may want a coarse-grained clustering of news stories into three or four clusters, while others may want a fine-grained clusterings into twenty or more. But in many cases, we will have choose K ourselves, with little outside guidance.

One solution is to choose the number of clusters to maximize some computable quantity of the clustering. First, note that the likelihood of the training data will always increase with K . For example, if a good solution is available for $K = 2$, then we can always obtain that same solution at $K > 2$; usually we can find an even better solution by fitting the data more closely. The Akaike Information

Criterion (AIC; Akaike, 1974) solves this problem by minimizing a linear combination of the log-likelihood and the number of model parameters, $AIC = 2m - 2\mathcal{L}$, where m is the number of parameters and \mathcal{L} is the log-likelihood. Since the number of parameters increases with the number of clusters K , the AIC may prefer more parsimonious models, even if they do not fit the data quite as well.

Another choice is to maximize the **predictive likelihood** on heldout data $\mathbf{x}_{1:N_h}^{(h)}$. This data is not used to estimate the model parameters ϕ and μ ; we can compute the predictive likelihood on this data by keeping the parameters ϕ and μ fixed, and running a single iteration of the E-step. In document clustering or **topic modeling** (Blei, 2012), a typical approach is to split each instance (document) in half. We use the first half to estimate $q_i(z_i)$, and then on the second half we compute the expected log-likelihood,

$$\ell_i = \sum_z q_i(z) (\log p(\mathbf{x}_i | z; \phi) + \log p(z; \mu)). \quad (4.24)$$

On heldout data, this quantity will not necessarily increase with the number of clusters K , because for high enough K , we are likely to overfit the training data. Thus, choosing K to maximize the predictive likelihood on heldout data will limit the extent of overfitting. Note that in general we cannot analytically find the K that maximizes either AIC or the predictive likelihood, so we must resort to grid search: trying a range of possible values of K , and choosing the best one.

Finally, it is worth mentioning an alternative approach, called **Bayesian nonparametrics**, in which the number of clusters K is treated as another latent variable. This enables statistical inference over a set of models with a variable number of clusters; this is not possible with EM, but there are several alternative inference procedures that are suitable for this case (Murphy, 2012), including MCMC (section 4.4). Reisinger and Mooney (2010) provide a nice example of Bayesian nonparametrics in NLP, applying it to unsupervised word sense induction.

4.3 Applications of EM

EM is not really an “algorithm” like, say, quicksort. Rather, it is a framework for learning with missing data. The recipe for using EM on a problem of interest is:

- Introduce latent variables \mathbf{z} , such that it is easy to write the probability $P(\mathcal{D}, \mathbf{z})$, where \mathcal{D} is your observed data; it should also be easy to estimate the associated parameters, given knowledge of \mathbf{z} .

- Derive the E-step updates for $q(\mathbf{z})$, which is typically factored as $q(\mathbf{z}) = \prod_i q_{z_i}(z_i)$, where i is an index over instances.
- The M-step updates typically correspond to the soft version of some supervised learning algorithm, like Naïve Bayes.

Some more applications of this basic setup are presented here.

Word sense clustering

In the “demos” folder, you can find a demonstration of expectation-maximization for word sense clustering. I assume we know that there are two senses, and that the senses can be distinguished by the contextual information in the document. The basic framework is identical to the clustering model of EM as presented above.

Semi-supervised learning

Nigam et al. (2000) offer another application of EM: **semi-supervised learning**. They apply this idea to document classification in the classic “20 Newsgroup” dataset, in which each document is a post from one of twenty newsgroups from the early days of the internet.

In the setting considered by Nigam et al. (2000), we have labels for some of the instances, $\langle \mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)} \rangle$, but not for others, $\langle \mathbf{x}^{(u)} \rangle$. The question they pose is: can unlabeled data improve learning? If so, then we might be able to get good performance from a smaller number of labeled instances, simply by incorporating a large number of unlabeled instances. This idea is called **semi-supervised learning**, because we are learning from a combination of labeled and unlabeled data; the setting is described in much more detail in chapter 20.

As in Naïve Bayes, the learning objective is to maximize the joint likelihood,

$$\log p(\mathbf{x}^{(\ell)}, \mathbf{x}^{(u)}, \mathbf{y}^{(\ell)}) = \log p(\mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)}) + \log p(\mathbf{x}^{(u)}) \quad (4.25)$$

We treat the labels of the unlabeled documents as missing data — in other words, as a latent variable. In the E-step we impute $q(y)$ for the unlabeled documents only. The M-step computes estimates of μ and ϕ from the sum of the observed counts from $\langle \mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)} \rangle$ and the expected counts from $\langle \mathbf{x}^{(u)} \rangle$ and $q(\mathbf{y})$.

Nigam et al. (2000) further parametrize this approach by weighting the unlabeled documents by a scalar λ , which is a tuning parameter. The resulting crite-

tion is:

$$\mathcal{L} = \log p(\mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)}) + \lambda \log p(\mathbf{x}^{(u)}) \quad (4.26)$$

$$\geq \log p(\mathbf{x}^{(\ell)}, \mathbf{y}^{(\ell)}) + \lambda E_q[\log p(\mathbf{x}^{(u)}, y)] \quad (4.27)$$

The scaling factor does not really have a probabilistic justification, but it can be important to getting good performance, especially when the amount of labeled data is small in comparison to the amount of unlabeled data. In that scenario, the risk is that the unlabeled data will dominate, causing the parameters to drift towards a “natural clustering” that may be a bad fit for the labeled data. Nigam et al. (2000) show that this approach can give substantial improvements in classification performance when the amount of labeled data is small.

Multi-component modeling

Now let us consider an alternative application of EM to supervised classification. One of the classes in 20 newsgroups is `comp.sys.mac.hardware`; suppose that within this newsgroup there are two kinds of posts: reviews of new hardware, and question-answer posts about hardware problems. The language in these **components** of the `mac.hardware` class might have little in common. So we might do better if we model these components separately. Nigam et al. (2000) show that EM can be applied to this setting as well.

Recall that Naïve Bayes is based on a generative process, which provides a stochastic explanation for the observed data. For multi-component modeling, we envision a slightly different generative process, incorporating both the observed label y_i and the latent component z_i :

- For each document i ,
 - draw the label $y_i \sim \text{Categorical}(\mu)$
 - draw the component $z_i \mid y_i \sim \text{Categorical}(\beta_{y_i})$, where $z_i \in 1, 2, \dots, K_z$.
 - draw the vector of counts $\mathbf{x}_i \mid z_i \sim \text{Multinomial}(\phi_{z_i})$

Our labeled data includes $\langle \mathbf{x}_i, y_i \rangle$, but not z_i , so this is another case of missing data. Again, we sum over the missing data, applying Jensen’s inequality to as to obtain a lower bound on the log-likelihood,

$$\log p(\mathbf{x}_i, y_i) = \log \sum_z^{K_z} p(\mathbf{x}_i, y_i, z) \quad (4.28)$$

$$\geq \log p(y_i; \mu) + E_q[\log p(\mathbf{x}_i \mid z; \phi) + \log p(z \mid y_i; \psi) - \log q_i(z)]. \quad (4.29)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

We are now ready to apply expectation-maximization. As usual, the distribution over the missing data — the component z_i — $q_i(z)$ is updated in the E-step. Then during the m-step, we compute:

$$\beta_{y,z} = \frac{E_q[\text{count}(y, z)]}{\sum_{z'} E_q[\text{count}(y, z')]} \quad (4.30)$$

$$\phi_{z,j} = \frac{E_q[\text{count}(z, j)]}{\sum_{j'} E_q[\text{count}(z, j')]} \quad (4.31)$$

Suppose we assume each class y is associated with K components, \mathcal{Z}_y . We can then add a constraint to the E-step so that $q_i(z) = 0$ if $z \notin \mathcal{Z}_y \wedge Y_i = y$.

4.4 Other approaches to learning with latent variables*

Expectation maximization is a very general way to think about learning with latent variables, but it has some limitations. One is the sensitivity to initialization, which means that we cannot simply run EM once and expect to get a good solution. Indeed, in practical applications of EM, quite a lot of attention may be devoted to finding a good initialization. A second issue is that EM tends to be easiest to apply in cases where the latent variables have a clear decomposition (in the cases we have considered, they decompose across the instances). For these reasons, it is worth briefly considering some alternatives to EM.

Sampling

Recall that in EM, we set $q(\mathbf{z}) = \prod_i q_i(z_i)$, factoring the q distribution into conditionally independent q_i distributions. In sampling-based algorithms, rather than maintaining a distribution over each latent variable, we draw random samples of the latent variables. If the sampling algorithm is designed correctly, this procedure will eventually converge to drawing samples from the true posterior, $p(\mathbf{z}_{1:N} \mid \mathbf{x}_{1:N})$. For example, in the case of clustering, we will draw samples from the distribution over clusterings of the data. If a single clustering is required, we can select the one with the highest joint likelihood, $p(\mathbf{z}_{1:N}, \mathbf{x}_{1:N})$.

This general family of algorithms is called **Markov Chain Monte Carlo** (MCMC): “Monte Carlo” because it is based on a series of random draws; “Markov Chain” because the sampling procedure must be designed such that each sample depends

only on the previous sample, and not on the entire sampling history. Gibbs Sampling is a particularly simple and effective MCMC algorithm, in which we sample each latent variable from its posterior distribution,

$$z_i \mid \mathbf{x}, \mathbf{z}_{-i} \sim p(z_i \mid \mathbf{x}, \mathbf{z}_{-i}), \quad (4.32)$$

where \mathbf{z}_{-i} indicates $\{\mathbf{z} \setminus z_i\}$, the set of all latent variables except for z_i .

What about the parameters, ϕ and μ ? One possibility is to turn them into latent variables too, by adding them to the generative story. This requires specifying a prior distribution; the Dirichlet is a typical choice of prior for the parameters of a multinomial, since it has support over vectors of non-negative numbers that sum to one, which is exactly the set of permissible parameters for a multinomial. For example,

$$\phi_y \sim \text{Dirichlet}(\alpha), \forall y \quad (4.33)$$

We can then sample $\phi_y \mid \mathbf{x}, \mathbf{z} \sim p(\phi_y \mid \mathbf{x}, \mathbf{z}, \alpha)$; this posterior distribution will also be Dirichlet, with parameters $\alpha + \sum_{i: y_i=y} \mathbf{x}_i$. Alternatively, we can analytically marginalize these parameters, as in **Collapsed Gibbs Sampling**; this is usually preferable if possible. Finally, we might maintain ϕ and μ as parameters rather than latent variables. We can employ sampling in the E-step of the EM algorithm, obtaining a hybrid algorithm called Monte Carlo Expectation Maximization (MCEM; Wei and Tanner, 1990).

In principle, these algorithms will eventually converge to the true posterior distribution. However, there is no way to know how long this will take; there is not even any way to check on whether the algorithm has converged. In practice, convergence again depends on initialization, since it might take ages to recover from a poor initialization. Thus, while Gibbs Sampling and other MCMC algorithms provide a powerful and flexible array of techniques for statistical inference in latent variable models, they are not a panacea for the problems experienced by EM.

Murphy (2012) includes an excellent chapter on MCMC; for a more comprehensive treatment, see Robert and Casella (2013).

Spectral learning

A more recent approach to learning with latent variables is based on the **method of moments**. In these approaches, we avoid the problem of non-convex log-likelihood by using a different estimation criterion. Let us write $\bar{\mathbf{x}}_i$ for the normalized vector of word counts in document i , so that $\bar{\mathbf{x}}_i = \mathbf{x}_i / \sum_j x_{ij}$. Then we

can form a matrix of word-word co-occurrence counts,

$$\mathbf{C} = \sum_i \mathbf{x}_i \mathbf{x}_i^\top. \quad (4.34)$$

We can also compute the expected value of this matrix under $p(\mathbf{x} \mid \phi, \mu)$, as

$$E[\mathbf{C}] = \sum_i \sum_k P(Z_i = k \mid \mu) \phi_k \phi_k^\top \quad (4.35)$$

$$= \sum_k N \mu_k \phi_k \phi_k^\top \quad (4.36)$$

$$= \Phi \text{Diag}(N\mu) \Phi^\top, \quad (4.37)$$

where Φ is formed by horizontally concatenating $\phi_1 \dots \phi_K$, and $\text{Diag}(N\mu)$ indicates a diagonal matrix with values $N\mu_k$ at position (k, k) . Now, by setting \mathbf{C} equal to its expectation, we obtain,

$$\mathbf{C} = \Phi \text{Diag}(N\mu) \Phi^\top, \quad (4.38)$$

which is very similar to the eigendecomposition $\mathbf{C} = \mathbf{Q} \mathbf{\Lambda} \mathbf{Q}^\top$. This suggests that simply by finding the eigenvectors and eigenvalues of \mathbf{C} , we could obtain the parameters ϕ and μ , and this is what motivates the name **spectral learning**.

However, there is a key difference in the constraints on the solutions to the two problems. In eigendecomposition, we require orthonormality, so that $\mathbf{Q} \mathbf{Q}^\top = \mathbf{I}$. But in estimating the parameters of a mixture model, we require the columns of Φ represents probability vectors, $\forall k, j, \phi_{k,j} \geq 0$, $\sum_j \phi_{k,j} = 1$, and that the entries of μ correspond to the probabilities over components. Thus, spectral learning algorithms must include a procedure for converting the solution into vectors of probabilities. One approach is to replace eigendecomposition (or the related singular value decomposition) with non-negative matrix factorization (Xu et al., 2003), which guarantees that the solutions are non-negative (Arora et al., 2013).

After obtaining the parameters ϕ and μ , we can obtain the distribution over clusters for each document by simply computing $p(z_i \mid \mathbf{x}_i; \phi, \mu) \propto p(\mathbf{x}_i \mid z_i; \phi) p(z_i; \mu)$. The advantages of spectral learning are that it obtains (provably) good solutions without regard to initialization, and that it can be quite fast in practice. Anandkumar et al. (2014) describe how similar matrix and tensor factorizations can be applied to statistical estimation in many other forms of latent variable models.

Chapter 5

Language models

A **language model** is used to compute the probability of a sequence of text. Why would we want to do this? Thus far, we have considered problems where text is the **input**, and we want to select an output, such as a document class or a word sense. But in many of the most prominent problems in language technology, text itself is the output:

- In **machine translation**, we convert text in a source language to text in a target language.
- In **speech recognition**, we convert audio signal to text.
- In **summarization**, we produce short texts that capture the key points of some longer text.

The goal of language models is to produce more **fluent** text output by computing the probability of the text.

Specifically, suppose we have a vocabulary of word types

$$\mathcal{V} = \{aardvark, abacus, \dots, zither\} \quad (5.1)$$

Given a sequence of word tokens w_1, w_2, \dots, w_M , with $w_i \in \mathcal{V}$, we would like to compute the probability $p(w_1, w_2, \dots, w_M)$. We will do this in a data-driven way, assuming we have a **corpus** of text.

- For now, we'll assume that the vocabulary \mathcal{V} covers all the word tokens that we will ever see. Of course, we can enforce this by allocating a special token $\langle \text{UNK} \rangle$ for unknown words. However, this might not be a great solution, as we will see later.

- Language models typically make an independence assumption across sentences, $p(s_1, s_2, \dots) = \prod_j p(s_j)$, where each sentence $s_j = [w_1, w_2, \dots, w_{M_j}]$. So for our purposes, it is sufficient to compute the probability of sentences. The justification for this assumption is that the probability of words that are not in the same sentence don't depend on each other too much. Clearly this isn't true: once I mention *Manuel Noriega* once in a document, I'm far more likely to mention him again (Church, 2000). But the dependencies between words within a sentence are usually even stronger, and are more relevant to the fluency considerations inherent in applications such as translation and speech recognition (which are typically evaluated at the sentence level anyway).

So how can we compute the probability of a sentence? The simplest idea would be to apply a **relative frequency estimator**:

$$p(\text{Computers are useless, they can only give you answers}) \quad (5.2)$$

$$= \frac{\text{count}(\text{Computers are useless, they can only give you answers})}{\text{count}(\text{all sentences ever spoken})} \quad (5.3)$$

It is useful to think about this estimator in terms of bias and variance.

- In the theoretical limit of infinite data, it might work. But in practice, we are asking for accurate counts over an infinite number of events, since sentences can be arbitrarily long.
- Even if we set an aggressive upper bound of, say, $n = 20$, the number of possible sentences is $\#|\mathcal{V}|^{20}$. A small vocabulary for English would have $\#|\mathcal{V}| = 10^4$, so we would have 10^{80} possible sentences.
- Clearly, this estimator is very data-hungry. We need to introduce bias to have a chance of making reliable estimates from finite training data.

Are language models meaningful? What are the probabilities of the following two sentences?

- *Colorless green ideas sleep furiously*
- *Furiously sleep ideas green colorless*

Noam Chomsky used this pair of examples to argue that the probability of a sentence is a meaningless concept:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Any English speaker can tell that the first sentence is grammatical but the second sentence is not.
- Yet neither sentence, nor their substrings, had ever appeared at the time that Chomsky wrote this article (they have appeared lots since then).
- Thus, he argued, empirical probabilities can't distinguish grammatical from ungrammatical sentences.

Pereira (2000) showed that by identifying *classes* of words (e.g., noun, verb, adjective, adverb — but not necessarily these grammatical categories), it is easy to show that the first sentence is more probable than the second. Class-based language models are discussed in section 5.4.

Are language models useful? Suppose we want to translate a sentence from Spanish:

(5.1) *El cafe negro me gusta mucho.*
The coffee black me pleases much.

But a good language model of English will tell us:

$$p(\textit{The coffee black me pleases much}) < p(\textit{I love dark coffee}) \quad (5.4)$$

How can we use this fact? Warren Weaver, one of the early leaders in machine translation, viewed it as a problem of breaking a secret code:

When I look at an article in Russian, I say: 'This is really written in English, but it has been coded in some strange symbols. I will now proceed to decode.'

This motivates a generative model (like Naïve Bayes!):

- English sentence $w^{(e)}$ generated from language model $p_{W^e}(w^{(e)})$
- Spanish sentence $w^{(s)}$ generated from noisy channel $p_{W^s|W^e}(w^{(s)} | w^{(e)})$

Then the **decoding** problem is:

$$\max_{w^{(e)}} p_{W^{(e)}|W^s}(w^{(e)} | w^{(s)}) \propto p_{W^s, W^e}(w^{(s)}, w^{(e)}) = p_{W^e}(w^{(e)}) p_{W^s|W^e}(w^{(s)} | w^{(e)}) \quad (5.5)$$

- The **translation model** is $p_{W^s|W^e}(w^{(s)} | w^{(e)})$. This ensures the **adequacy** of the translation.

- The **language model** is $p_{W^e}(w^{(e)})$. This ensures the **fluency** of the translation.

What else can we model with a noisy channel?

- Speech recognition (original = words; encoded = sound)
- Spelling correction (original = well-spelled text; encoded = text with spelling mistakes)
- Part of speech tagging (original = tags; encoded = words)
- Parsing (original = parse tree; encoded = words)
- ...

The key insight of the noisy channel model is that it allows us to decompose NLP systems into a translation model and a language model. Since the language model be estimated from unlabeled data, this means we can improve our system without the expense of obtaining more labeled data — we simply focus on improving $p_{W^e}(w)$. It also means that language models are in principle reusable across many language technology systems. For this reason, I will focus on language models in this chapter, and return to machine translation later in the course.

5.1 N-gram language models

Let us return to the relative frequency estimator,

$$p(\text{Computers are useless, they can only give you answers}) \quad (5.6)$$

$$= \frac{\text{count}(\text{Computers are useless, they can only give you answers})}{\text{count}(\text{all sentences ever spoken})} \quad (5.7)$$

We'll define the probability of a sentence as the probability of the words (in order): $p(w) = p(w_1, w_2, \dots, w_M)$. We can apply the chain rule:

$$\begin{aligned} p(w) &= p(w_1, w_2, \dots, w_M) \\ &= p(w_1)p(w_2 \mid w_1)p(w_3 \mid w_2, w_1) \dots p(w_M \mid w_{M-1}, \dots, w_1) \end{aligned}$$

Each element in the product is the probability of a word given all its predecessors. We can think of this as a *word prediction* task: *Computers are [BLANK]*. The

relative frequency estimate of the probability of the word *useless* in this context is,

$$p(\textit{useless} \mid \textit{computers are}) = \frac{\text{count}(\textit{computers are useless})}{\sum_x \text{count}(\textit{computers are } x)} = \frac{\text{count}(\textit{computers are useless})}{\text{count}(\textit{computers are})}.$$

Note that we haven't made any approximations yet, and we could have just as well applied the chain rule in reverse order, $p(\mathbf{w}) = p(w_M)p(w_{M-1} \mid w_M) \dots$, or in any other order. But this means that we also haven't really improved anything either: to compute the conditional probability $P(w_M \mid w_{M-1}, w_{M-2}, \dots)$, we need to model V^{M-1} contexts, with V possible events. We can't even **store** this probability distribution, let alone reliably estimate it.

N-gram models make a simple approximation: condition on only the past $n-1$ words.

$$p(w_m \mid w_{m-1} \dots w_1) \approx P(w_m \mid w_{m-1}, \dots, w_{m-n+1}) \quad (5.8)$$

This means that the probability of a sentence \mathbf{w} can be computed as

$$p(w_1, \dots, w_M) \approx \prod_m p(w_m \mid w_{m-1}, \dots, w_{m-n+1}) \quad (5.9)$$

To compute the probability of an entire sentence, it is convenient to pad the beginning and end with special symbols $\langle \text{START} \rangle$ and $\langle \text{STOP} \rangle$. Then the bigram ($n = 2$) approximation to the probability of *I like black coffee* is:

$$p(I \mid \langle \text{START} \rangle) p(\textit{like} \mid I) p(\textit{black} \mid \textit{like}) p(\textit{coffee} \mid \textit{black}) p(\langle \text{STOP} \rangle \mid \textit{coffee}) \quad (5.10)$$

In this model, we have to estimate and store the probability of only V^n events (exponential in the order of the n-gram), and not V^M (exponential in the length of the sentence).

The n-gram probabilities can be determined by relative frequency estimation,

$$Pr(W_i = c \mid W_{i-1} = b, W_{i-2} = a) = \frac{\text{count}(a, b, c)}{\sum_{c'} \text{count}(a, b, c')} = \frac{\text{count}(a, b, c)}{\text{count}(a, b)} \quad (5.11)$$

In estimation, there could be at two problems with an n -gram language model:

- **n is too small.** In this case, we are missing important linguistic context. Consider the following sentences:

(5.2) *Gorillas always like to groom **THEIR** friends.*

(c) Jacob Eisenstein 2014-2015. Work in progress.

(5.3) The **computer** that's on the 3rd floor of our office building **CRASHED**.

The uppercase bolded words depend crucially on their predecessors in lowercase bold: the likelihood of *their* depends on knowing that *gorillas* is plural, and the likelihood of *crashed* depends on knowing that the subject is a *computer*. If the n -grams are not big enough to capture this context, then the resulting language model would offer probabilities that are too low for these sentences, and too high for sentences that fail basic linguistic tests like number agreement.

- **n is too big.** In this case, we can't make good estimates of the n -gram parameters from our dataset, because of data sparsity. To handle the *gorilla* example, we would need to model 6-grams; which means accounting for V^6 events. Under a very small vocabulary of $V = 10^4$, this means estimating the probability of 10^{24} distinct events.

These two problems point to another **bias/variance** tradeoff. Can you see how it works? In practice, we often have **both** problems at the same time. Language is full of long-range dependencies that we cannot capture because n is too small; at the same time, language datasets are full of rare phenomena, whose probabilities we fail to estimate accurately because n is too large.

We will seek approaches to keep n large, while still making low-variance estimates of the underlying parameters. To do this, we will introduce a different sort of bias: **smoothing**. But before we talk about that, let's consider how we can evaluate language models.

5.2 Evaluating language models

Because language models are typically components of larger systems — language modeling is not really an application itself — we would prefer **extrinsic evaluation**. This means evaluating whether the language model improves performance on the application task, such as machine translation or speech recognition. But this is often hard to do, and depends on details of the overall system which may be irrelevant to language modeling. In contrast, **intrinsic evaluation** is task-neutral. Better performance on intrinsic metrics may be expected to improve extrinsic metrics across a variety of tasks, unless we are over-optimizing the intrinsic metric. We will discuss intrinsic metrics here, but bear in mind that it is important to also perform extrinsic evaluations to ensure that the improvements obtained on these intrinsic metrics really carry over to the applications that we care about.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Held-out likelihood

A popular intrinsic metric is the **held-out likelihood**.

- We obtain a test corpus, and compute the (log) probability according to our model. As in classification, it is crucial that the sentences in the held-out corpus were not used in estimating the model itself.
- A good language model should assign high probability to this held-out data. Specifically, we compute,

$$\ell(\mathbf{w}) = \sum_i^N \sum_m^{M_i} \log p(w_m^{(i)} | w_{m-1}^{(i)}, \dots, w_{m-n+1}^{(i)}), \quad (5.12)$$

summing over all sentences $\{\mathbf{w}^{(i)}\}_{i \in 1 \dots N}$ in the held-out corpus.

Typically, unknown words in the test data are mapped to the $\langle \text{UNK} \rangle$ token. This means that we have to estimate some probability for $\langle \text{UNK} \rangle$ on the training data. One way to do this is to fix the vocabulary \mathcal{V} to the $V - 1$ words with the highest counts in the training data, and then convert all other tokens to $\langle \text{UNK} \rangle$.

Perplexity

Perplexity is a transformation of the held-out likelihood into an information-theoretic quantity. Specifically, we compute

$$PP(\mathbf{w}) = 2^{-\frac{\ell(\mathbf{w})}{M}}, \quad (5.13)$$

where M is the total number of tokens in the held-out corpus.

- The transformation means that lower perplexities correspond to higher likelihoods, so lower scores are better on this metric. (Lower perplexity is better, because you are less perplexed.) In the limit, we obtain probability 1 for our held-out corpus, with $PP = 2^{-\log 1} = 2^0 = 1$.
- Assume a uniform, unigram model in which $p(w_i) = \frac{1}{V}$ for all V words in the vocabulary. Then,

$$\begin{aligned} PP(\mathbf{w}) &= \left[\left(\frac{1}{V} \right)^M \right]^{-\frac{1}{M}} \\ &= \left(\frac{1}{V} \right)^{-1} = V \end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

These observations imply that we can think of perplexity as the *weighted branching factor* at each word in the sentence.

- If we have solved the word prediction problem perfectly, $PP(\mathbf{w}) = 1$, because there is only one possible choice for each word.
- If we have a uniform model that assigns equal probability to every word, then $PP(\mathbf{w}) = V$. This is not a worst-case scenario — in the worst case, we assign **zero** probability to some word in the test data — but it is a worst “reasonable” case.
- Most models give perplexities that fall somewhere in between 1 and V .

Example On 38M tokens of WSJ, $V \approx 20K$, (Jurafsky and Martin, 2009, page 97) obtain these perplexities on a 1.5M token test set.

- Unigram: 962
- Bigram: 170
- Trigram: 109

Will it keep going down?

Information theory*

Perplexity is closely related to the concept of **entropy**, the expected value of the information contained in each word.

$$H(P) = - \sum_w p(\mathbf{w}) \log p(\mathbf{w}) \quad (5.14)$$

The true entropy of English (or any real language) is unknown. Claude Shannon, one of the founders of information theory, wanted to compute upper and lower bounds. He would read passages of 15 characters to his wife, and ask her to guess the next character, recording the number of guesses it took for her to get the correct answer. As a fluent speaker of English, his wife could provide a reasonably tight bound on the number of guesses needed per character.¹

We can view the goal of language modeling as computing a distribution Q that is similar to the true distribution P . To measure the quality of Q , we can compute

¹Question for you: is this an upper bound or a lower bound?

its **cross-entropy** with P , written as $H(P, Q)$,

$$H(P, Q) = E_P[\log Q] \quad (5.15)$$

$$= - \sum_{\mathbf{w}} p(\mathbf{w}) \log q(\mathbf{w}) \quad (5.16)$$

$$= - \sum_{\mathbf{w}} p(\mathbf{w}) \log \left(q(\mathbf{w}) \frac{p(\mathbf{w})}{p(\mathbf{w})} \right) \quad (5.17)$$

$$= - \sum_{\mathbf{w}} p(\mathbf{w}) \log \frac{q(\mathbf{w})}{p(\mathbf{w})} + p(\mathbf{w}) \log p(\mathbf{w}) \quad (5.18)$$

$$= \sum_{\mathbf{w}} p(\mathbf{w}) \log \frac{p(\mathbf{w})}{q(\mathbf{w})} - p(\mathbf{w}) \log p(\mathbf{w}) \quad (5.19)$$

$$= D_{KL}(P||Q) + H(P), \quad (5.20)$$

where $D_{KL}(P||Q)$ is the **Kullback-Leibler (KL) divergence** between P and Q . The KL-divergence is a non-symmetric measure of the the dissimilarity of two distributions, where $\forall(P, Q), D_{KL}(P||Q) \geq 0$ and $D_{KL}(P||P) = 0$.² The cross-entropy also includes a term for the entropy of the true distribution P , but since P is given, we can only control Q . Thus, minimizing the cross entropy $H(P, Q)$ is equivalent to minimizing the KL-divergence $D_{KL}(P||Q)$.

We do not have access to the true P , just a sequence $\mathbf{w} = \{w_1, w_2, \dots\}$, which is sampled from P . In the limit, the length of \mathbf{w} is infinite, so we have,

$$H(P, Q) = - \sum_{\mathbf{w}} p(\mathbf{w}) \log q(\mathbf{w}) \quad (5.21)$$

$$= - \lim_{M \rightarrow \infty} \frac{1}{M} \log q(\mathbf{w}). \quad (5.22)$$

There term $p(\mathbf{w})$ disappears because the word sequence \mathbf{w} is itself a sample from this distribution. In practice, we have finite M , so we compute the approximation,

$$H(P, Q) \approx - \frac{1}{M} \log q(\mathbf{w}) \quad (5.23)$$

$$PP(Q) = 2^{-\frac{1}{M} \log q(\mathbf{w})} = 2^{H(P, Q)} \quad (5.24)$$

Thus, the perplexity of the language model Q can be derived from its cross-entropy with the true word distribution P , which we estimate from the observed

²KL-divergence has connections to expectation maximization: the lower bound on the expected likelihood can be viewed as the true likelihood minus the KL-divergence $D_{KL}(q(\mathbf{y})||p(\mathbf{y} | \mathbf{x}))$, so that the E-step minimizes the KL-divergence by setting $q(\mathbf{y}) = p(\mathbf{y} | \mathbf{x})$.

word sequence w . Low perplexity implies low cross-entropy, which in turn implies a low KL-divergence between P and Q .

Further aside A related topic in psycholinguistics is the “constant entropy rate hypothesis,” also called the “uniform information density hypothesis.” The hypothesis is that speakers should prefer linguistic choices that convey a uniform amount of information over time (Jaeger, 2010). Some evidence:

- Speakers shorten predictable words, and lengthen unpredictable ones (Jaeger, 2010).
- Low-probability words slow down the reader (Smith and Levy, 2013)
- Syntactic reductions (e.g., *I’m* versus *I am*) are more likely when the reducible word contains less information (Jaeger and Levy, 2006).

5.3 Smoothing and discounting

Limited data is a persistent problem in estimating language models. In section 5.1, we presented n-grams as a partial solution. But as we saw, sparse data can be a problem even for low-order n-grams; at the same time, many linguistic phenomena, like subject-verb agreement, cannot be incorporated into language models without higher-order n-grams. It is therefore necessary to add additional inductive biases to n-gram language models. This section covers some of the most intuitive and common approaches, but there are many more. Chen and Goodman (1999) provides a good survey of the state-of-the-art in the late 1990s; more recent approaches are discussed in section 5.4.

Smoothing

A major concern in language modeling is to avoid the situation $p(w) = 0$, which could arise as a result of a single unseen n-gram. A similar problem arose in Naïve Bayes, and there we solved it by **smoothing**: adding pseudo counts. The same idea can be applied to n-gram language models, as shown here in the bigram case,

$$p_{\text{Laplace}}(W_i = b \mid W_i = a) = \frac{\text{count}(a, b) + \alpha}{\sum_{w'} \text{count}(a, w') + V\alpha}. \quad (5.25)$$

- In general, this is called **Lidstone smoothing**.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- When $\alpha = 1$, it is **Laplace smoothing**.
- When $\alpha = 0.5$, we are following **Jeffreys-Perks law**.
- Manning and Schütze (1999) offer more insight on the justifications for Jeffreys-Perks smoothing

To maintain normalization, anything that we add to the numerator (α) must also appear in the dominator ($V\alpha$). This idea is reflected in the concept of **effective counts**:

$$c_i^* = (c_i + \alpha) \frac{N}{N + V\alpha}, \quad (5.26)$$

where c_i is the count of event i , and c_i^* is the effective count. The **discount** for each n-gram is then computed as,

$$d_i = \frac{c_i^*}{c_i} = \frac{(c_i + \alpha)}{c_i} \frac{N}{(N + \alpha)}$$

Discounting and backoff

Discounting “borrows” probability mass from observed n-grams and redistributes it. In Lidstone smoothing, we borrow probability mass by increasing the denominator of the relative frequency estimates, and redistribute it by increasing the numerator for all n-grams. But instead, we could borrow the same amount of probability mass from all observed counts, and redistribute it among only the unobserved counts. This is called **absolute discounting**.

For example, if we set an absolute discount $d = 0.1$ in a trigram model, we get: $p(w \mid \text{denied the}) =$

word	counts c	effective counts c^*	unsmoothed probability	smoothed probability
<i>allegations</i>	3	2.9	0.429	0.414
<i>reports</i>	2	1.9	0.286	0.271
<i>claims</i>	1	0.9	0.143	0.129
<i>request</i>	1	0.9	0.143	0.129
<i>charges</i>	0	0.2	0.000	0.029
<i>benefits</i>	0	0.2	0.000	0.029
...				

We need not redistribute the probability mass equally. Instead, we can **back-off** to a lower-order language model. In other words: if you have trigrams, use

trigrams; if you don't have trigrams, use bigrams; if you don't even have bigrams, use unigrams. This is called **Katz backoff**:

$$c^*(u, v) = c(u, v) - d \quad (5.27)$$

$$p_{\text{backoff}}(v | u) = \begin{cases} \frac{c^*(u, v)}{c(u)} & \text{if } c(u, v) > 0 \\ \alpha(u) \times \frac{p_{\text{backoff}}(v)}{\sum_{v': c(u, v')=0} p_{\text{backoff}}(v')} & \text{if } c(u, v) = 0 \end{cases} \quad (5.28)$$

Typically we can set the discount d to minimize perplexity on a development set.

Interpolation

An alternative to discounting is **interpolation**: setting the probability of a word in context to a weighted sum of its probabilities across progressively shorter contexts.

Instead of choosing a single n for the size of the n -gram, we can take the weighted average across several n -gram probabilities,

$$\begin{aligned} p_{\text{Interpolation}}(a | b, c) &= \lambda_3^{(a)} p_3^*(a | b, c) \\ &\quad + \lambda_2^{(a)} p_2^*(a | b) \\ &\quad + \lambda_1^{(a)} p_1^*(a). \end{aligned}$$

In this equation, p_k^* is the maximum likelihood estimate (MLE) of a k -gram model, and $\lambda_k^{(a)}$ is the weight of the n -gram model p_k^* for word a . A nice property of this model is that it can learn to use longer context for some words (e.g., possessive pronouns like *his* and *her*, which often match the gender of the entity as defined earlier in the sentence), and shorter context for others (e.g., rare content words).

To ensure that the interpolated $p(w)$ is still a probability, we have a constraint, $\sum_k \lambda_k^{(a)} = 1, \forall a$. But how to find the specific values of λ for each word? An elegant solution is **expectation maximization**. Recall from chapter 4 that we can think about EM as learning with **missing data**: we just need to choose missing data such that learning would be easy if it weren't missing. So what's missing in this case? We can think of each word w_m as drawn from an n -gram of unknown size, $z_m \in \{1 \dots n\}$. This z_m is the missing data that we are looking for! Specifically, consider the following generative story:

- For each token m ,

(c) Jacob Eisenstein 2014-2015. Work in progress.

- draw $z_m \sim \text{Categorical}(\lambda^{(w_m)})$
- draw $w_m \sim p_{z_m}^*(w_m \mid w_{m-1}, \dots, w_{m-z_m})$.

If we knew $\{z_m\}_{m \in 1 \dots M}$, then we could compute λ from relative frequency estimation, $\lambda_k^{(a)} = \frac{\sum_m \delta(z_m=k) \delta(w_m=a)}{\sum_m \delta(w_m=a)}$.³ Since we do not know the values of the missing data, we impute a distribution $q_m(z_m)$ in the E-step, which represents our degree of belief that word token w_m was generated from a n-gram of order z_m .

Having defined these quantities, we can derive EM updates:

• **E-step:**

$$q_m(k) = \text{Pr}(Z_m = k \mid \mathbf{w}_{1:m}) \quad (5.29)$$

$$\propto p_z^*(w_m \mid w_{m-1}, \dots, w_{m-k+1}) \lambda_k^{(w_m)} \quad (5.30)$$

• **M-step:**

$$\lambda_k(a) = \frac{E_q[\text{count}(W = a, Z = k)]}{\sum_{k'} E_q[\text{count}(W = a, Z = k')]} \quad (5.31)$$

$$= \frac{\sum_m q_m(k) \delta(w_m = a)}{\sum_m \delta(w_m = a)} \quad (5.32)$$

As usual, EM iterates between these two steps until convergence to a local optimum.

Kneser-Ney smoothing

Kneser-Ney smoothing also incorporates discounting, but redistributes the resulting probability mass in a different way. Consider the example: *I recently visited*

- *Francisco?*
- *Duluth?*

Now suppose that both bigrams *visited Duluth* and *visited Francisco* are unobserved in our training data, and furthermore, that the unigram probability $p^*(\text{Francisco})$ is greater than $p^*(\text{Duluth})$. Nonetheless we would still guess that $p(\text{visited Duluth}) > p(\text{visited Francisco})$, because *Duluth* is a more **versatile** word. This notion of versatility is the key to Kneser-Ney smoothing.

³We could also use z to update our n-gram models p_i^* , but we will assume those are fixed here.

Writing c for a context of undefined length, and $\text{count}(w, c)$ as the count of word w in context c , we define the Kneser-Ney bigram probability as

$$p_{KN}(w | c) = \begin{cases} \frac{\text{count}(w, c) - d}{\text{count}(c)}, & \text{count}(w, c) > 0 \\ \alpha(c) p_{\text{continuation}}(w), & \text{otherwise} \end{cases}$$

$$p_{\text{continuation}}(w) = \frac{\#|c : \text{count}(w, c) > 0|}{\sum_{w'} \#|c' : \text{count}(c', w') > 0|}$$

First, note that we reserve probability mass using absolute discounting d , which is taken from all unobserved n-grams. The total amount of discounting in context c is $d \times \#|w : \text{count}(w, c) > 0|$, and we divide this equally among the unseen n-grams,

$$\alpha(c) = \frac{d \times \#|w : \text{count}(w, c) > 0|}{\#|c : \text{count}(w, c) = 0|}. \quad (5.33)$$

This is the amount of probability mass left to account for versatility, which we define via the *continuation probability* $p_{\text{continuation}}(w)$ as proportional to the number of observed contexts in which w appears. In the numerator of the continuation probability we have the number of contexts c in which w appears, and in the denominator, we normalize by computing the same quantity over all words w' .

In practice, interpolation works a little better than backoff,

$$p_{KN}(w | c) = \frac{\text{count}(w, c) - d}{\text{count}(c)} + \lambda_c p_{\text{continuation}}(w) \quad (5.34)$$

This idea of counting contexts may seem heuristic, but there is a cool theoretical justification from Bayesian nonparametrics (Teh, 2006).

5.4 Other Types of N-gram Language Models

Interpolated Kneser-Ney is pretty close to state-of-the-art. But there are some interesting other types of language models, and they apply ideas that we have already learned.

Class-based language models

The reason we need smoothing is because even the trigram probability model $P(w_m, | w_{m-1}, w_{m-2})$ has a huge number of parameters. We could use the idea of

(c) Jacob Eisenstein 2014-2015. Work in progress.

word classes to simplify. Imagine that each word has a latent class z ,

$$p_{\text{class}}(w_m \mid w_{m-1}) = \sum_z p(w_m \mid z; \phi) p(z \mid w_{m-1}; \beta),$$

where $z \in [1, K]$, $K \ll V$. This means that each word w_m is conditioned on its class z through parameter ϕ_z , and the class itself is conditioned on the previous word w_{m-1} through parameter $\beta_{w_{m-1}}$. The advantage of this approach is that it gives a bigram probability using $2 \times V \times K$ parameters, instead of V^2 .

How do we estimate such a model? Since there is missing data — the word classes — we might use expectation maximization:

- E-step: update $q_i(z)$
- M-step: update ϕ and β

But this is usually too slow in practice, since it requires multiple passes over the training data, which is typically very large. A useful approximate algorithm is **exchange clustering** (Brown et al., 1992), which assigns each word type to a single class, rather than maintaining a soft distribution $p(z \mid w; \beta)$. This algorithm incrementally constructs a binary tree over the words in the vocabulary, so that each word can be represented by a bit vector corresponding to the series of left/right decisions to get to the word from the root. The prefixes of these bit vectors are an early form of **word embedding**, and it has been shown that syntactically similar words tend to have similar bit vectors, as shown in Figure 5.1. As we will see in chapter 20, these vectors can be used as features in NLP systems, improving their performance by enabling generalization from frequent to rare words (Miller et al., 2004; Koo et al., 2008).

Discriminative language models

Alternatively, we could just train a model to predict $p(w_m \mid w_{m-1}, w_{m-2}, \dots)$ directly. We can think of this as a straightforward classification problem, where the label space is equal to the entire vocabulary; for example, Rosenfeld (1996) applies logistic regression to language modeling, and Roark et al. (2007) apply perceptrons and conditional random fields (section 9.4). A key advantage is that discriminative training minimizes the error rate, rather than maximizing probability; for applications such as speech recognition, this is a better fit for the ultimate goal, which is recognizing speech with as few errors as possible. Moreover, because the underlying model is now discriminative, additional features can be

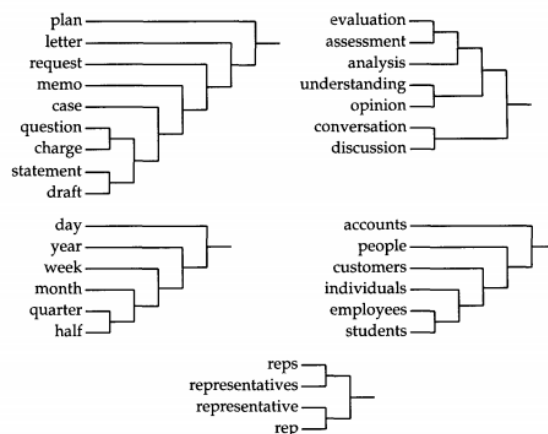


Figure 5.1: Example subtrees from the Brown et al. (1992) hierarchical class-based language model

included, such as features of the syntactic structure (Khudanpur and Wu, 2000). However, Roark et al. (2007) report that discriminative language models are expensive to train, requiring extensive feature selection.

Neural language models

Currently, **neural probabilistic language models** are attracting a lot of interest. These are related to discriminative language models, but they also maintain a continuous state that can capture long-term history. Another key distinction in these models is the use of dense, discriminatively-trained vector representations, computing the probability $p(w_m \mid \mathbf{w}_{1:m-1})$ by passing an inner product $\mathbf{h}_m^\top \mathbf{v}_{w_{m-1}}$ through a sigmoid activation function.⁴

There are many flavors of neural language models, from the early formulation by Bengio et al. (2003) to more elaborate contemporary models based on long short-term memory (LSTM; Hochreiter and Schmidhuber, 1997; Sundermeyer et al., 2012). We will focus on the Recurrent Neural Network Language Model (RNNLM; Mikolov et al., 2010), which works fairly well in practice, and is the basis for some of the more complex recent models.

⁴The function $p_{ij} = \frac{\exp \psi_{ij}}{\sum_{j'} \exp \psi_{ij'}}$ is sometimes called **softmax** in the neural net literature.

Assume each word type i is associated with a dense vector representation \mathbf{u}_i , which is a parameter of the model. Writing \mathbf{x}_m is an indicator vector such that

$$x_{mi} = \begin{cases} 1, & i = w_m \\ 0, & i \neq w_m, \end{cases} \quad (5.35)$$

then we can write $\mathbf{u}_{w_m} = \mathbf{U}\mathbf{x}_m$. In fact, we will have two dense vector representations per word: \mathbf{U} for the input and \mathbf{V} for the output. We will return to word vectors in chapter 15 and chapter 21.

The RNNLM is **recurrent** in the sense that there is a hidden state \mathbf{h}_m at each word position m , which is constructed from the hidden state \mathbf{h}_{m-1} , as well as from the word w_m . Specifically, the hidden state in the RNNLM is given by,

$$\mathbf{h}_m = f(\mathbf{U}\mathbf{x}_m + \Theta\mathbf{h}_{m-1}), \quad (5.36)$$

where f is an element-wise non-linear activation function, such as the sigmoid. Finally, we predict w_{m+1} with probability

$$p(w_{m+1} \mid \mathbf{h}_m) = \frac{\exp(\mathbf{h}_m^\top \mathbf{v}_{w_{m+1}})}{\sum_i \exp(\mathbf{h}_m^\top \mathbf{v}_i)}. \quad (5.37)$$

Since the hidden states $\mathbf{h}_{1:m}$ can be computed deterministically from the words $\mathbf{w}_{1:m}$, the RNNLM defines a distribution $p(w_{m+1} \mid \mathbf{w}_{1:m})$ without any explicit limit on the length of the past history. However, information from words $w_j : j \ll m$ will be attenuated by repeatedly passing through the recurrent function. Recent variants on the RNNLM address this issue through the use of memory cells (Sundermeyer et al., 2012) and gates (Chung et al., 2015), enabling crucial pieces of past information to more directly impact future predictions.

The RNNLM has three parameters: the word representation matrices \mathbf{U} and \mathbf{V} , and the recurrent update matrix Θ . Note that the size of these parameters are relatively small. Writing K for the size of each word vector representation (assuming the input and output representations have identical size), the parameter dimensions are $K \times V$ for \mathbf{U} and \mathbf{V} , and $K \times K$ for Θ . Considering that $V \gg K$ in all practical cases, this means that the RNNLM has far fewer parameters than even a bigram language model, which must store counts of size $V \times V$. For this reason, neural language models require less attention to smoothing and regularization than traditional n-gram language models.

The RNNLM parameters are learned by **backpropagation** from a loss function: a typical choice is the negative log-likelihood of the data, which is identical to the

cross-entropy:

$$\ell = \sum_m \log p(w_m \mid \mathbf{w}_{1:m-1}) \quad (5.38)$$

$$= \sum_m \mathbf{h}_{m-1}^\top \mathbf{v}_{w_m} - \log \sum_i^V \exp(\mathbf{h}_{m-1}^\top \mathbf{v}_i). \quad (5.39)$$

Computing this loss function (and its gradients) can be expensive, since it involves summing over the entire vocabulary at each word position. One alternative is to use a **hierarchical softmax** function to compute the sum more efficiently, in $\log V$ time (Mikolov et al., 2011); another is to optimize an alternative metric, such as noise-contrastive estimation, which learns by distinguishing observed instances from artificial instances generated from a noise distribution (Mnih and Teh, 2012).

Recent work on **probabilistic programming** has resulted in a number of toolkits for building computation graphs over architectures such as the RNNLM. These toolkits — Theano and Torch are currently popular choices — perform automatic differentiation, allowing the researcher to plug in a variety of different loss functions and model architectures, without having to derive and implement the parameter updates by hand.

5.5 Other details

Datasets Dataset genre is important: an LM learned from Shakespeare is a poor match for the Wall Street Journal (WSJ); an LM learned from the WSJ is a poor choice for predictive text entry in cellphones.

Vocabulary We have assumed we know the total vocabulary size V . Will we always know this? What if we don't?

- Suppose we are estimating a bigram language model. Then we can set $V_{\text{bigram}} = V_{\text{unigram}}^2$, assuming we have seen all unigrams.
- But in general, we're always at risk of seeing new words (<http://www.americandialect.org/hashtag-2012>)
 - *hashtag* (2012 word of the year)
 - *phablet*

(c) Jacob Eisenstein 2014-2015. Work in progress.

- *Gangnam*⁵
- -3.78109932019384

- If the set of unigrams is defined in advance, this is the **closed vocabulary** setting. Typically we will just replace unknown words with a special token, $\langle \text{unk} \rangle$.
- Another solution is to backoff from the unigram model to a character model:

$$P_u(s_i) = \begin{cases} \frac{\text{count}(s_i)}{\text{count}(\text{all tokens}) + \beta}, & s_i \in \mathcal{V} \\ \beta P_c(s_i), & s_i \notin \mathcal{V} \end{cases} \quad (5.40)$$

$$P_c(s_i) = P_{\text{len}}(s_i) \prod_{a_j \in s_i} \frac{\text{count}(a_j)}{\text{count}(\text{all characters})} \quad (5.41)$$

[**todo: reconcile this notation with the rest of the chapter**] We could even have a bigram or trigram model over characters.

- Still another possibility, this time in the setting of neural language models, would be to require that word representations are themselves composed from a character-level RNNLM (Ling et al., 2015). Besides ensuring that we can always compute a word representation for any sequence of character symbols, this approach implies that similarly-spelled words have similar representations.

⁵[**todo: Wow this list is already incredibly dated! I should probably try to use more timeless examples, like *Sputnik* or something.**]

Part II

Sequences and trees

Chapter 6

Morphology

So far we have been focusing on NLP at the word level. Now we will explore meaning **inside of words**. We've already hinted at a morphological problem by introducing the idea of **lemmas**, where *serve/served/serving* all have the lemma *serve*.

From the perspective of document classification, these multiple forms may just seem like an annoyance, which we can get rid of by lemmatization or stemming (more on this later). But morphology conveys information which can be crucial for some applications.

Information retrieval With a search query like *bagel*, we want to get hits for the **inflected** form *bagels*; the same goes for irregular inflections like *corpus/corpora*, *goose/geese*. In **query expansion**, the search query is expanded to include all inflections of the search terms. Note that this isn't always what we want: for example, given a query for *Apple*, we may not want hits for *apples*.

Information extraction A major goal of information extraction is to capture references to events, and their properties. Event timing is conveyed in morphology: in English, we have suffixes for past tense (*she talked*), the past participle (*she had spoken*), and the present participle (*she is speaking*). Other languages can indicate many more details about event timing through morphology; for example, Romance languages like French have a much larger inventory of verb endings:

<i>J'achete un velo</i>	I buy a bicycle (now)
<i>J'acheterai un velo</i>	I will buy a bicycle
<i>J'achetais un velo</i>	I was buying a bicycle
<i>J'ai acheté un velo</i>	I bought a bicycle
<i>J'acheterais un velo</i>	I would buy a bicycle

In English, this function is mostly filled by auxiliary verbs like *will*, *was*, *had*, and *would*. This makes morphological analysis relatively less important for English, as we can get a long way with carefully constructed n-gram patterns (Riloff, 1996). But in languages like French and Spanish — where second-language learners are tormented by conjugation tables with dozens of different inflections — there seems little alternative to morphological analysis if language technology is to generalize across many verbs.

Document classification Even document classification tasks, such as sentiment analysis, are potentially impacted by morphology. For example, suppose you are doing sentiment analysis, and you encounter the out-of-vocabulary words *unfriended*, *antichrist*, *unputdownable*, or *disenchanted*. As unknown words, they would make no contribution to the overall sentiment polarity in a bag-of-words system. But with some morphological reasoning, we can see that they are indeed strongly subjective.

Translation In addition to recognizing morphology, there are applications in which we need to produce it. Translation is a classic case, especially when translating from morphologically simple languages like English and Chinese to morphologically rich languages, like French, Czech, German, and Swahili. Here again, a purely word-based approach would suffer from data sparsity: relatively rare words would be unlikely to be seen in every inflection, and thus the translation system would be unable to produce them.

Morphology, Orthography, and Phonology

Morphology interacts closely with two related systems: orthography and phonology. The **surface form** of a word is the form that is written down or spoken. This form results from the interactions between morphology and the orthographic and phonological systems. More specifically:

- **Morphology** describes how meaning is constructed from combining affixes. For example, it is a morphological fact of English that adding the affix +S to

Surface form	lemma	features
<i>duck</i>	<i>duck</i>	NOUN+SINGULAR
<i>ducks</i>	<i>duck</i>	NOUN+PLURAL
<i>duck</i>	<i>duck</i>	VERB+PRESENT
<i>ducks</i>	<i>duck</i>	VERB+THIRDPERSON+PRESENT

Table 6.1: Fragment of a morphologically-aware dictionary

many nouns creates a plural.

$$\textit{berry} + \text{PLURAL} \rightarrow \textit{berry+s}$$

Morphological rules may also include stem changes, such as *goose*+PLURAL \rightarrow *geese*.

- **Orthography** specifically relates to writing. For example,

$$\textit{berry+s} \rightarrow \textit{berries}$$

is an orthographic rule. We have lots of these in English, which is one reason English spelling is difficult.

- **Phonology** describes how sounds combine. For example, the different pronunciations of the final *s* in *cats* (s) and *dogs* (z) follow from a phonological rule (Bender, 2013, example 25, page 30).

In English, morphologically distinct words may be pronounced differently even when they are spelled the same, and this can reflect morphological differences. *read*+PRESENT vs. *read*+PAST. Conversely, morphological variants may be spelled differently even when they sound the same, like *The Champions' league* versus *The Champion's league* versus *The Champions league*.

Productivity

One idea for dealing with morphology is to build a morphologically-aware dictionary. The keys in this dictionary would correspond to **surface forms**, such as *served*. The values would include both the underlying **lemma** as well as any morphological features: in this case, the lemma is **serve**, and the feature is PAST. Given such a dictionary, we simply look up each surface form that we encounter.

As shown in the example in Table 6.1, we may need multiple entries for the same surface form; this means that there is ambiguity, so simple lookup will not suffice. Still another problem is that morphology is **productive**, meaning that it applies to new words. If you only know the words *Google* or *iPad*, you can immediately understand their inflected forms.

- Have you Googled that yet?
- I have broken all three iPads.

Derivational morphology (more on this later) is productive in another way: you can produce new words by applying morphological changes to existing words. hyper+un+desire+able+ity

In some languages, derivational morphology can create extremely complicated words. Jurafsky and Martin (2009) have a fun example from Turkish:

A Turkish word

uygarlaştıramadıklarımızdanmışsınızcasına

uygar_laş_tır_ama_dık_lar_ımız_dan_mış_sınız_casına

“as if you are among those whom we were not able to civilize (=cause to become civilized)”

uygar: *civilized*

_laş: *become*

_tır: *cause somebody to do something*

_ama: *not able*

_dık: *past participle*

_lar: *plural*

_ımız: *1st person plural possessive (our)*

_dan: *among (ablative case)*

_mış: *past*

_sınız: *2nd person plural (you)*

K. Oflazer pc to J&M

Figure 6.1: From (Jurafsky and Martin, 2009)

In the homework, you’ll see examples from Swahili, which also has complex morphology. A dictionary of all possible surface forms in such languages would be gargantuan. So instead of building a static dictionary, we will try to model the underlying morphological and orthographic rules.

(c) Jacob Eisenstein 2014-2015. Work in progress.

6.1 Types of morphemes

There are two broad classes of morphemes: **stems** and **affixes**. Intuitively, stems are the “main” part of meaning, and affixes are the modifiers. Typically, **stems** can appear on their own (they are **free**) and affixes cannot (they are **bound**).

Affixes can be categorized by where they appear with respect to the stem.

- **Prefixes:** *un+learn*, *pre+view*.
 - These examples are **derivational**, in that they form new words, rather than forming grammatical variants of the same word (*inflectional* morphology; more on this in section 6.2).
 - English has no inflectional prefixes, but other languages do. For example, in Swahili, *u-na-kata* means *you are cutting*, while *u-me-kata* means *you have cut*. In this example, *na* and *me* are prefixes, *kata* is the root.¹
- **Suffixes** are the typical way of inflecting words in English, and in other languages in the Indo-European family. For example, in English: *I learn+ed*, *She learn+s*, *three apple+s*, *four fox+es*. English suffixes can also be derivational: for example: *modern+ity*, *fix+able*, and *deriv+ation+al*.
- **Circumfixes** go around the stem.
 - German has a circumfix for the past participle: *sagen* (**say**) → *ge+sag+t* (**said**)
 - English has a very small number of circumfix examples: *bold* → *em+bold+en*, and, arguably, *light* → *en+light+en*. Both of these examples are derivational.
 - French negation can be seen as a circumfix: *Je mange+NEG* → *Je ne mange pas* (I do not eat).²
 - More generally, morphemes can be non-contiguous, e.g. (Bender, 2013, example 7, page 12):

¹Would it be better to think about *u*, *na*, and *me* as words? This example suggests that the word/affix distinction is not always clear-cut.

²In spoken French, the *ne* is gradually disappearing, so that *Je mange pas* is now acceptable.

(7)	Root	Pattern	Part of Speech	Phonological Form	Orthographic Form	Gloss
	ktb	CaCaC	(v)	katav	כתב	'wrote'
	ktb	hiCCiC	(v)	hixtiv	הכתוב	'dictated'
	ktb	miCCaC	(n)	mixtav	מכתב	'a letter'
	ktb	CCaC	(n)	ktav	כתב	'writing, alphabet'

[heb]

In this example, the root *ktb* (related to writing) is combined with patterns that indicate where to insert vowels to produce different parts-of-speech and meanings.

- **Infixes** go inside the stem.
 - In Tagalog (spoken in the Philippines), the root *hingi* indicates a request, and the infix *um* creates *humingi*, as in *I asked*.
 - English, *absolutely+fucking*→

(6.1) *absofuckinglutely*

(6.2) *?absfuckingsolutely*

where the '?' prefix indicates questionable linguistic acceptability.

- Morphology may be **non-segmental**, meaning that it doesn't involve any affix at all. For example, the pluralization of *goose* to *geese* is not accomplished through any affix, but through vowel alteration; the past tense marking of *eat* → *ate* is another example of this phenomenon, known as *apophony*. Languages in which morphemes are represented by affixes that are "glued together" (like *talk+ed* or *think+ing*) are known as **agglutinative**; languages in which morphemes are represented by changes to spelling and sound are known as **fusional**.
- What about words like *fish*, which have the same form in both singular and plural? We say that this word has a **zero** plural.

6.2 Types of morphology

Morphology serves a variety of linguistic functions, and acts in a variety of ways. Inflectional and derivational morphology are distinguished by their function; other

(c) Jacob Eisenstein 2014-2015. Work in progress.

forms of morphology, such as cliticization and compounding are distinguished by how they work. In this section, we will focus mainly on inflectional and derivational morphology, describing their roles in English, and in other languages when there is no adequate example in English.

Inflectional morphology

Inflectional morphology adds information about the stem, typically grammatical properties such as tense, number, and case. English has a relatively simple system of inflectional morphology, compared to many other languages.

Affix	Syntactic/semantic effect	Examples
-s	NUMBER: plural	<i>cats</i>
-'s	possessive	<i>cat's</i>
-s	TENSE: present, SUBJ: 3sg	<i>jumps</i>
-ed	TENSE: past	<i>jumped</i>
-ed/-en	ASPECT: perfective	<i>eaten</i>
-ing	ASPECT: progressive	<i>jumping</i>
-er	comparative	<i>smaller</i>
-est	superlative	<i>smallest</i>

Figure 6.2: From (Bender, 2013)

Nouns

English nouns are marked for **number** and **possession**. Number is typically marked by the suffix +s, e.g., *hat* + PLURAL → *hat+s*, but some words are pluralized differently, e.g., *geese*, *children*, and *fish*. Number is binary in English (singular versus plural), but many languages, such as Arabic and Sanskrit, include an additional **dual** number for groups of two. English has residual traces of the dual number, with *both* versus *all* and *either* versus *any*. Some Austronesian languages even have a **trial** number, for groups of three, and languages such as Arabic have a **paucal** number, for small groups. Conversely, nouns are not marked for number at all in Japanese and Indonesian.

Many languages mark nouns for **case**, which is the syntactic role that the noun plays in the sentence. In English, we do distinguish the case of some pronouns:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- *He* (NOMINATIVE) *gave her* (OBLIQUE) *his* (GENITIVE) *guitar*.
- *She gave him her guitar*.
- *I gave you our guitar*.
- *You gave me your guitar*.

The third person masculine pronoun appears as *he* in the nominative case, *him* in the oblique case, and *his* in the genitive case. English distinguishes these cases for all personal pronouns except for the second person, where the nominative and oblique cases are both *you*.

Other languages — such as Latin, Russian, Sanskrit, and Tamil — mark the case of all nouns. These languages have additional cases, such as dative (indirect object), accusative (direct object), and vocative (address). In German, noun is not inflected for case, but the articles and adjectives are, as shown in example 49 from Bender (2013):

- (6.3) *Der alte Mann gab dem kleinen Affen die grosse Banane.*
 The old man (NOM) gave the little monkey (DATIVE) the big banana (AC-
 CUSATIVE)

Notice how *der*, *dem*, and *die* all mean the same thing (*the*), but they are spelled differently due to the case marking. The adjectives (*alte*, *kleinen*, *grosse*) are also marked for case.

Many languages — such as Romance languages — mark the gender and number of nouns by inflecting the article and adjective. e.g., Spanish:

- (6.4) *El coche rojo pasó la luz roja.*
 The red car ran the red light.
- (6.5) *Los coches rojos pasó las luces rojas.*
 The red cars ran the red lights.

Here, *la* is the feminine article and *el* is the masculine article; the adjective for *red* is inflected to *roja* when describing a feminine noun (*luz*, meaning light), and *rojo* when describing a masculine noun (*coche*, meaning car). The article and adjective must **agree** with noun for the sentence to be grammatical. The following examples are ungrammatical for this reason:

- (6.6) **Los coches rojo pasó la luce rojas*
- (6.7) **Los coches rojas pasó las luces rojos*

(c) Jacob Eisenstein 2014-2015. Work in progress.

In English, demonstrative determiners mark number: e.g., *this book* vs *these books*, and the determiner and noun must agree, e.g. **this books*. Agreement is also required between subject and verb, as we will see shortly.

Romance languages like Spanish and French mark gender as masculine and feminine, but it need not be binary:

- English pronouns include neuter *it*; German, Sanskrit, and Latin do this for all nouns.
- Danish and Dutch distinguish **neuter** from **common** gender.[**todo: example**]
- Other languages distinguish **animate** and **inanimate** genders.

Verbs

English verbs are inflected for tense and number distinguishing past (*she acted*), present (*you act*), and third person singular (*she acts*). As with nouns, these inflections may change the orthography (*plan+ed* → *planned*), and there are many irregular patterns, e.g. *they eat* / *she eats* / *we ate*. English verbs are also inflected for aspect, distinguishing the perfective (*I had eaten*) and progressive (*I am eating*). The perfective and the past tense are identical for regular verbs, e.g. *we had talked*, *we talked*.

Many languages (e.g., Chinese and Indonesian), do not mark tense with morphology. For example, Indonesian uses function words rather than morphology to distinguish tense (Table 6.2).

<i>Saya makan apel</i>	I eat an apple
<i>Saya sedang makan apel</i>	I am eating an apple
<i>Saya telah makan apel</i>	I already ate an apple
<i>Saya akan makan apel</i>	I will eat an apple

Table 6.2: Indonesian uses function words (*sedang*, *telah*, *makan*) rather than morphology to distinguish verb tense. [**todo: switch to exe**]

Romance languages distinguish many more tenses than English with morphology. For example, Spanish has multiple past tenses: **preterite** and **imperfect**, distinguishing events that occurred at a specific past point in time from a continuous or repeated past state:

- (6.8) *I ate onions yesterday*
Comí cebollas ayer

- (6.9) *I ate onions every day*
Comía cebollas cada día

Spanish and French also have endings for conditional (*comería cebollas*, *I would eat onions*) and future (*comeré cebollas*, *I will eat onions*). In English, these differences are marked with time signals rather than morphology. In French and Spanish, time signals are also an option, e.g. *voy a comer cebollas*, which literally translates to *I am going to eat onions*.

Romance languages also have separate verb forms for every combination of number and person, while in English, only the third-person singular is distinguished:

- English: *I speak / you speak / she speaks / we speak / you (pl) speak / they speak*
- Spanish: *Yo hablo / tu hablas / ella habla / nosotros hablamos / vosotros hablais / ellas hablan*
- French: *Je parle / tu parles / elle parle / nous parlons / vous parlez / ils parlent*

In Spanish and in many other Romance languages (but not French), the verb morphology is sufficiently descriptive that the subject is often omitted, since it can often be easily recovered from the verb ending and the context.

Other things can be marked with affixes, such as **evidentiality** – how the speaker came to know the information. In Eastern Pomo (a California language), there are verb suffixes for four evidential categories (McLendon, 2003):

-ink'e	nonvisual sensory
-ine	inferential
-le	hearsay
-ya	direct knowledge

Adjectives and adverbs

Adjectives in English mark comparative and superlative (*taller*, *tallest*). Adverbs can mark comparative and superlative too: *Yangfeng paddles fast*, *Yi paddles faster*, *Uma paddles fastest*. As we have seen, adjectives can mark gender and number in languages like French and Spanish, where they are required to agree with the noun and determiner; adjectives also mark case in languages like German and Latin.

Synthetic and isolating languages

Languages with complex morphology are called **synthetic**; languages with simple morphology are called **isolating** or **analytic**. The **index of synthesis** quantifies this property by measuring the ratio of the number of morphemes in a given text to the number of words. On this index, English is relatively, but not extremely, analytic.

Language	Index of synthesis
Vietnamese	1.06
Yoruba	1.09
English	1.68
Old English	2.12
Swahili	2.55
Turkish	2.86
Russian	3.33
Inuit (Eskimo)	3.72

Figure 6.3: From Bender (2013)

An approximation of the index of synthesis is the type-token ratio. Can you see why? If you count the number of unique surface forms in 10K *parallel* sentences from a corpus of European Parliament transcripts, you get:

- English: 16k distinct word types
- French: 22k
- German: 32k
- Finnish: 55k

Derivational Morphology

Derivational morphology is a way to create new words and change part-of-speech.

- **nominalization**
 - *V + -ation: computerization*
 - *V + -er: walker*

(c) Jacob Eisenstein 2014-2015. Work in progress.

- *Adj + -ness: fussiness*
- *Adj + -ity: obesity*
- **negation:** *undo, unseen, misnomer*
- **adjectivization:** *V + -able : doable, thinkable, N + -al : tonal, national, N + -ous: famous, glamorous*
- **abverbization:** *ADJ + -ily: clumsily*
- **lots more:** *rewrite, phallocentrism, ...*

You can create totally new words this way.

word → *wordify* → *wordification* → *wordificationism* → *antiwordificationism* → *hyperantiwordificationism*

As with inflection, derivational morphology can require orthographic changes, e.g. *true+ly* → *truly* and *fussy+ness* → *fussiness*. It can also cause phonological changes, such as the change emphasis from *imPOSSible* to *impossiBILity*, and the change in vowel from *ferTILE* to *ferTILity*.

Other types of morphology

Cliticization combines *Georgia+’s* into *Georgia’s*; the possessive clitic *’s* is syntactically independent but phonologically dependent. This syntactic independence can be seen in examples like (Bender, 2013, example 21):

(6.10) Jesse met the president of the university’s cousin

In this example, the possessive modifies the *president*, but it attaches to the right edge of the entire noun phrase.

- Pronouns appear as clitics in French, e.g., *j’accuse* (I accuse), as does negation *Je n’accuse personne* (I don’t accuse anyone).
- Another example is from Hebrew: *l’shana tova* (literally for year good, meaning happy new year); the preposition *for* appears as a clitic.

Compounding combines two words into a new word:

(6.11) *cream* → *ice cream*

We can think of *ice cream* as a word since it is a non-compositional combination of *ice* and *cream*. Perhaps someday the written space will be dropped, as it has been in *watermelon* (Figure 6.4).

(c) Jacob Eisenstein 2014-2015. Work in progress.

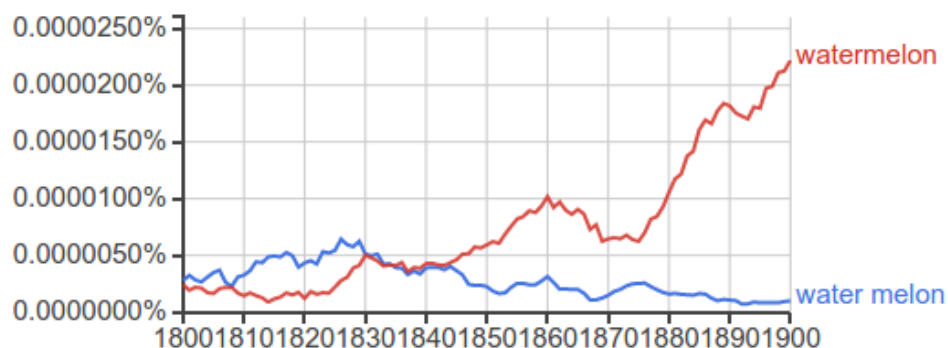


Figure 6.4: The written space in watermelon disappeared as the word became more frequent over the 19th century. From Google ngrams.

Portmanteaus combine words, truncating one or both.

(6.12) *smoke* + *fog* → *smog*

(6.13) *glass* + *asshole* → *glasshole*

Urban Dictionary is a fun source of contemporary portmanteaus.

Irregularities

English morphology contains a lot of irregularities: *know/knew/known*, *foot/feet*, *go/went*, etc. if you are not a native speaker, learning these was probably a pain in the neck. The good news is that there are fewer of these all the time! English is undergoing a process in which these irregular forms are gradually being replaced: for example, the past tense of *show* used to be *shew*, just as the past tense of *know* is still *knew* (Figure 6.5a). This transformation remains incomplete, as the past participle of *show* is still *shown*, and not *showed* (Figure 6.5b). However, this example points to the bad news for language learners: the most frequently-occurring words, like *know*, will be the last to change — if ever!

6.3 Computing and morphology

In this section, we will briefly overview some of the computational problems related to morphology. We don't yet have many tools to solve these problems, but

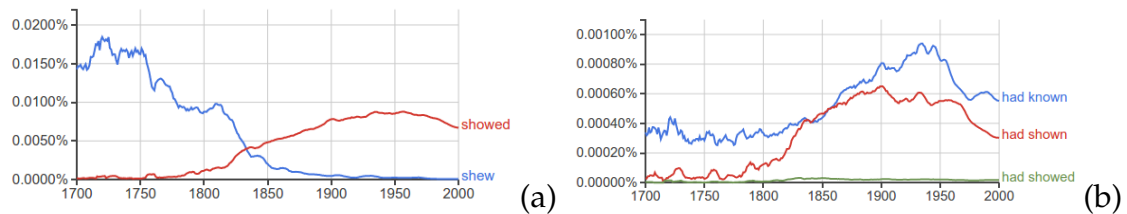


Figure 6.5: Google n-grams plots for inflections of *show*. While the past participle *had shown* is decreasing, this does not seem to be due to competition from the more regular *had showed*; rather, there appears to be a broader decrease in frequency of the past participle, shown by the parallel pattern for *had known*.

we will soon: chapter 7 presents finite-state automata, which are the workhorse of morphological analysis in NLP. For now, we will simply state the problem definitions, and discuss some of the challenges involved.

Lemmatization

[todo: write]

Stemming

[todo: write]

Generation

[todo: write]

Normalization

[todo: write]

(c) Jacob Eisenstein 2014-2015. Work in progress.

Chapter 7

Finite-state automata

Consider the following problems:

- Segment a word into its stem and affixes: *impossibility* \rightarrow *im+possibl+ity*.
- Convert a sequence of morphemes like *im+possible+ity* into the correct sequence of characters (*impossibility*).
- Decide whether a given word is morphotactically correct, or more generally, rank all the possible realizations for a morphological expression like NEGATION + *possible*: *impossible, inpossible, nonpossible, unpossible*, etc.
- Given a speech utterance and a large set of potential text transcriptions, choose the one with the highest probability according to an n-gram language model.
- Perform context-sensitive spelling correction, so as to correct examples like *their at piece* to *they're at peace*.

All of these problems relate to the content of the previous two chapters — language models and morphology — but none of them seem easily solved by supervised classifiers. This chapter presents a new tool for language technology: finite state automata. Finite-state automata are particularly suited for scoring strings (sequences of characters, words, morphemes, or phonemes), and for converting one string into another. A key advantage of finite state automata is their modularity: the output of one finite-state transducer can be the input for another, allowing the combination of simple components into cascades with rich and complex behaviors.

Finite-state automata are a formalism for representing a subset of formal languages, the **regular** languages; these are languages that can be defined with regu-

lar expressions. While there is strong evidence that natural language is not regular — that is, the question of whether a given sentence is grammatical cannot be answered with any regular expression — finite state automata can be used as the building block for a surprisingly wide range of applications in language technology.¹

7.1 Automata and languages

Finite state automata emerge from formal language theory. Here are some basic formalisms that will be used throughout this chapter:

- An **alphabet** Σ is a set of symbols, e.g. $\{a, b, c, \dots, z\}$, or $\{aardvark, abacus, \dots, zyxt\}$.
- A **string** ω is a sequence of symbols, $\omega \in \Sigma^*$. The empty string ϵ contains zero symbols.
- A **language** $L \subseteq \Sigma^*$ is a set of strings.
- An **automaton** is an abstract model of a computer, which reads a string $\omega \in \Sigma^*$, and determines whether or not $\omega \in L$.

This seems a very different notion of “language” than English or Hindi. But could we think of these natural languages in the same way as formal languages? If *impossible* is acceptable as an English word but *unpossible* is not, might it be possible to build an automaton that formalizes the underlying linguistic distinction?

Finite-state automata

A finite-state **acceptor** (FSA) is a special type of automaton, which is capable of modeling some, but not all languages. Formally, finite-state automata are defined by a tuple $M = \langle Q, \Sigma, q_0, F, \delta \rangle$, consisting of:

- a finite alphabet Σ of input symbols;
- a finite set of **states** $Q = \{q_0, q_1, \dots, q_n\}$;

¹A more formal treatment of finite state automata and their applications to language is offered by Mohri et al. (2002). Knight and May (2009) show how finite-state automata can be composed together to create impressive applications, focusing on **transliteration** of words and names between languages with different scripts. Here, we’ll build the formalism from the ground up, starting with finite-state acceptors, then adding weights, and then adding transduction, finally arriving at the same sorts of applications.

- a **start state** $q_0 \in Q$;
- a set of **final states** $F \subseteq Q$;
- a **transition function** $\delta : Q \times \Sigma \rightarrow 2^Q$. The transition function maps from a state and an input symbol to a **set** of possible resulting states.

Given this definition, M accepts a string ω if there is a path from q_0 to any state $q_i \in F$ that consumes all of the symbols in ω . If M accepts ω , this means that ω is in the formal language L defined by M .

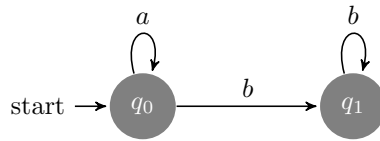
Example Consider the following FSA, M_1 .

$$\Sigma = \{a, b\} \quad (7.1)$$

$$Q = \{q_0, q_1\} \quad (7.2)$$

$$F = \{q_1\} \quad (7.3)$$

$$\delta = \{ \{ (q_0, a) \rightarrow q_0 \}, \{ (q_0, b) \rightarrow q_1 \}, \{ (q_1, b) \rightarrow q_1 \} \} \quad (7.4)$$



This FSA defines a language over an alphabet of two symbols, a and b . The transition function δ is written as a set of tuples: the tuple $\{(q_0, a) \rightarrow q_0\}$ says that if you are in state q_0 and you see symbol a , you can consume it and stay in q_0 . Because each pair of initial state and symbol has at most one resulting state, this FSA is **deterministic**: each string ω induces at most one path. Note that δ does not contain any information about what to do if you encounter the symbol a while in state q_1 . In this case, you are stuck, and cannot accept the input string.

What strings does this FSA accept? We begin in q_0 , but we have to get to q_1 , since this is the only final state. We can accept any number of a symbols while in q_0 , but we require a b symbol to transition to q_1 . Once there, we can accept any number of b symbols, but if we see an a symbol, there is nothing we can do. So the regular expression corresponding to the language defined by M_1 is a^*bb^* . To see this, consider what M_1 would do if it were fed each of the following strings: $aaabb$; aa ; $abbb$; bb .

Regular languages* Can every formal language be recognized by some finite state automata? No. Finite state automata can only recognize **regular languages**. The classic example of a non-regular language is $a^n b^n$; this language includes only those strings that contain n copies of symbol a , followed by n copies of symbol b . The **pumping lemma** demonstrates that this language cannot be accepted by any FSA. The proof is by contradiction. Suppose M is an FSA that accepts the language $a^n b^n$. By definition M must have a finite number of states; if we choose a string $a^m b^m$ such that m is bigger than the number of states in M , then the path through M must contain a cycle, and the transitions on this cycle must accept only the symbol a . But if there is a cycle, then we can repeat the cycle any number of times, “pumping up” the number of a symbols in the string. The automaton M must therefore also accept strings $a^{m'} b^m$, with $m' > m$. But these strings are not in the language $a^n b^n$, so we arrive at a contradiction. The proof will be covered in detail by any textbook on theory of computation (e.g., Sipser, 2012).

Determinism

- In a deterministic (D)FSA, the transition function is defined so that $\delta : Q \times \Sigma \rightarrow Q$. This means that every pair of initial state and symbol can transition to at most one resulting state.
- In a nondeterministic (N)FSA, $\delta : Q \times \Sigma \rightarrow 2^Q$. This means that a pair of initial state and symbol can transition to multiple resulting states. As a consequence, an NFSA may have multiple paths to accept a given string.
- We can determinize any NFSA using the powerset construction, but the number of states in the resulting DFSA may be exponential in the size of the original NFSA.
- Any **regular expression** can be converted into an NFSA, and thus into a DFSA.

The English Dictionary as an FSA We can build a simple “chain” FSA which accepts any single word. So, we can define the English dictionary with an FSA. However, we can make this FSA much more compact. (see slides)

- Begin by taking the **union** of all of the chain FSAs by defining **epsilon transitions** (transitions that do not consume an input symbol) from the start state to chain FSAs for each word (5303 states / 5302 arcs using a 850 word dictionary of “basic English”).

- Eliminate the epsilon transitions by pushing the first letter to the front (4454 states / 4453 arcs)
- **Determinize** (2609 / 2608)
- **Minimize** (744 / 1535). The cost of minimizing an acyclic FSA is $O(E)$. This data structure is called a **trie**.

Operations In discussing talked about three operations: union, determinization and minimization. Other important operations are:

Intersection only accept strings in both FSAs: $\omega \in (M_1 \cap M_2)$ iff $\omega \in M_1 \cap \omega \in M_2$.

Negation only accept strings not accepted by FSA M : $\omega \in (\neg M)$ iff $\omega \notin M$.

concatenation accept strings of the form $\omega = [\omega_1\omega_2]$, where $\omega_1 \in M_1$ and $\omega_2 \in M_2$.

FSAs are **closed** under all these operations, meaning that resulting automaton is still an FSA (and therefore still defines a regular language).

FSAs for Morphology

Now for some morphology. Suppose that we want to write a program that accepts only those words that are constructed in accordance with English derivational morphology:

- *grace, graceful, gracefully*
- *disgrace, disgraceful, disgracefully, ...*
- *Google, Googler, Googleology, ...*
- **gracelyful, *disungracefully, ...*

As we saw in the English dictionary example, we could just make a list, and then take the union of the list using ϵ -transitions. The list would get very long, and it would not account for productivity (our ability to make new words like *antiwordificationist*). So let's try to use finite state machines instead. Our FSA will have to encode rules about morpheme ordering, called *morphotactics*.

Every word must have a stem, so we do not want to accept proposed words like *dis-* or *-ly*. This suggests that we should have at least two states: one for before we have seen a stem, and one for after. Assuming the alphabet Σ consists of all English morphemes, we can define a transition function so that it is only possible to transition from q_0 to q_1 by consuming a stem morpheme; by defining

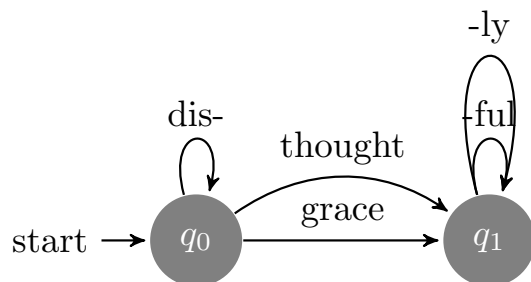


Figure 7.1: First try at modeling English morphology

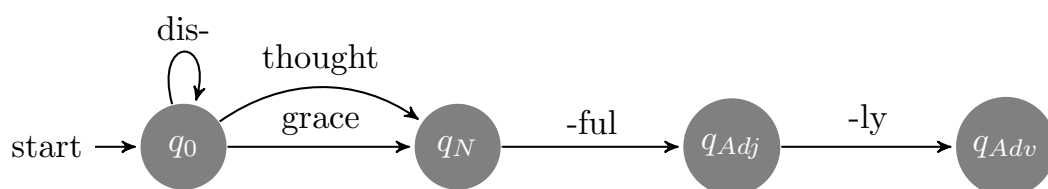


Figure 7.2: Second try at modeling English morphology, this time distinguishing parts-of-speech

If $F = \{q_1\}$, we can ensure that every word has a stem. For prefixes, we can allow self-transitions in q_0 on prefix morphemes; we can do the same in q_1 for suffix morphemes.

The resulting FSA is shown in Figure 7.1 will accept *grace*, *disgrace*, *graceful*, *disgraceful*, and even *disgracefully* (with two self-transitions in q_1). However, it will also accept **gracelyful* and **gracerly*. To deal with these cases, we need to think about what the suffixes are doing. The suffix *-ful* converts the noun *grace* into an adjective *graceful*; it does the same for words like *thoughtful* and *sinful*. The suffix *-ly* converts the adjective *graceful* to the adverb *gracefully* (to see the difference, compare *the ballet was graceful* to *the ballerina moved gracefully*.) These examples suggest that we need additional states in our FSA, such as q_{noun} , $q_{\text{adjective}}$, and q_{adverb} . Each of these is a potential final state, and the suffixes allow transitions between them. This FSA is shown in Figure 7.2.

However, with a little more thought, we see that this approach is still too simple. First, not every noun can be made into an adjective: **chairful* and **monkeyful* are perhaps suggestive of some kind of poetic meaning, but would not be recognized as standard English. Second, many nouns are made into adjectives using different suffixes, such as *music+al*, *fish+y*, and *elv+ish*. We need to create ad-

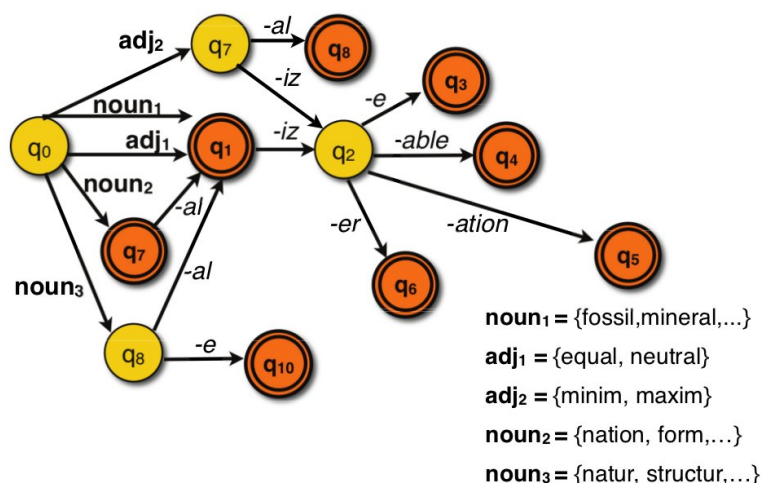


Figure 7.3: A fragment of a finite-state acceptor for derivational morphology. From Julia Hockenmaier's slides.

ditional noun states to distinguish these noun groups, so as to avoid accepting ill-formed words like **musicky* and **fishful*. We could continue to refine the FSA, coming ever-closer to an accurate model of English morphotactics. A fragment of such an FSA is shown in Figure 7.3.

This approach makes a key assumption: every word is either in or out of the language, with no wiggle room. Perhaps you agreed that *musicky* and *fishful* were not valid English words; but if forced to choose, you probably find *a fishful stew* or *a musicky tribute* preferable to *behaving disgracefully*. To take the argument further, here are some Google counts for various derivational forms:

- *superfast*: 70M; *ultrafast*: 16M; *hyperfast*: 350K; *megafast*: 87K
- *suckitude*: 426K; *suckiness*: 378K
- *nonobvious*: 1.1M; *unobvious*: 826K; *disobvious*: 5K

Given this diversity of possible realizations of the same idea, rather than asking whether a word is **acceptable**, we might like to ask how acceptable it is. But finite state acceptors gives us no way to express *preferences* among technically valid choices. We will need to augment the formalism for this.

7.2 Weighted Finite State Automata

A weighted finite-state automaton $M = \langle Q, \Sigma, \pi, \xi, \delta \rangle$ consists of:

- A finite set of states $Q = \{q_0, q_1, \dots, q_n\}$
- A finite alphabet Σ of input symbols
- Initial weight function, $\pi : Q \rightarrow \mathbb{R}$
- Final weight function $\xi : Q \rightarrow \mathbb{R}$
- A transition function $\delta : Q \times \Sigma \times Q \rightarrow \mathbb{R}$

We have departed from the FSA formalism in three ways:

- Every state can be a start state, with score π_q .
- Every state can be an end state, with score ξ_q .
- Transitions are possible between any pair of states on any input, with a score $\delta_{q_i, \omega, q_j}$.

Now, we can score every path through a weighted finite state acceptor (WFSA) by the sum of the weights for the transitions, plus the scores for the initial and final states. The **shortest path algorithm** finds the minimum-cost path through a WFSA for a string ω , with time complexity $\mathcal{O}(E + V \log V)$, where E is the number of edges and V is the number of vertices (Cormen et al., 2009).

Weighted finite state automata (WFSAs) are a generalization of unweighted FSAs: for any FSA M we can build an equivalent WFSa by setting $\pi_q = \infty$ for all $q \neq q_0$, $\xi_q = \infty$ for all $q \notin F$, and $\delta_{q_i, \omega, q_j}$ for all transitions $\{(q_1, \omega) \rightarrow q_2\}$ that are not permitted by the transition function of M .

Applications of WFSAs

We can use WFSAs to score derivational morphology as suggested above. But let's start with some simpler examples.

Edit distance

An **edit distance** is a function of two strings, which quantifies their similarity: for example, *she* and *he* differ by only the addition of a single letter, while *you* and *me* differ on every letter. There are a huge number of ways to compute edit

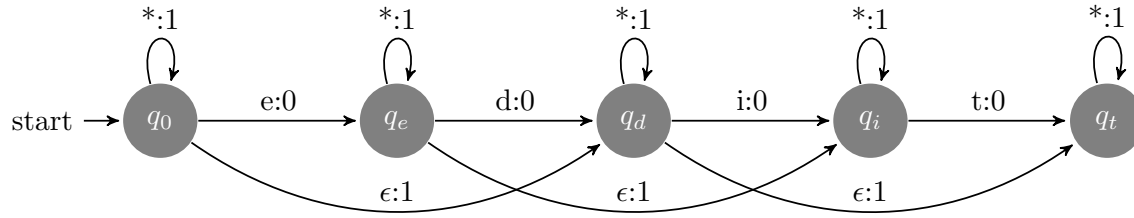


Figure 7.4: A weighted finite state acceptor for computing edit distance from the word *edit*.

distance (Manning et al., 2008), with applications in information retrieval, bioinformatics, and beyond.

Here we consider a simple edit distance, which computes the minimum number of character insertions, deletions, and substitutions required to get from one word to another. Insertions and deletions are penalized by a cost of one; substitutions have a cost of two. To compute this cost, we build a WFA with one state for every letter in the word, plus an initial state q_0 : for example, for the word *edit*, we build a machine with states q_0, q_e, q_d, q_i, q_t .

- The initial cost for q_0 is zero; for every other state, the initial cost is infinite.
- The final cost for q_t is zero; for every other state, the final cost is infinite.
- We define the transition function as follows:
 - The cost for “correct” symbols and rightward moves is zero: for example, $\delta_{q_0,e,q_e} = 0$, and $\delta_{q_i,t,q_t} = 0$.
 - The cost for self-transitions is one, regardless of the symbol: for example, $\delta_{q_d,*,q_d} = 1$. These self-transitions correspond to **insertions**.
 - The cost for epsilon transitions to the right is one: for example, $\delta_{q_e,\epsilon,q_d} = 1$. These transitions correspond to **deletions**.
 - The cost of all other transitions is ∞ .

The machine is shown in Figure 7.4. The total edit distance for a string is the *sum* of costs across the best path through machine. Note that we did not define a cost for **substitutions** (e.g., from *him* to *ham*), because substitutions can be performed by a combination of insertion and deletion, for a total cost of two. However, some edit distances assign a cost of one to substitutions; can you see how to modify the WFA to compute such an edit distance?

N-gram language models

Weighted finite state acceptors can also be used to compute probabilities of sequences — for example, the probability of a word sequence from an n-gram language model. To do this, we define the states and transitions so that each transition is equal to a condition probability, $\delta_{q_i, \omega_m, q_j} = p(q_i, \omega_m \mid q_j)$, so that the product is equal to the joint probability of the state sequence and the string,

$$p(\mathbf{q}_{1:M}, \boldsymbol{\omega}_{1:M}) = \prod_m^M p(q_m, \omega_m \mid q_{m-1}). \quad (7.5)$$

For example, to construct a unigram language model over a vocabulary \mathcal{V} of size V , we need just a single state. All transitions are self-transitions, with probability equal to the unigram word probability, $\delta_{q_0, w, q_0} = p_1(w)$.

To construct a bigram language model, we need to model the conditional probability $p(w_m \mid w_{m-1})$. To do this in a WFSA, we must create V different states: one for each context. Then we define the transition function as,

$$\delta_{q_i, w_m, q_j} = \begin{cases} p(w_m \mid w_{m-1} = i), & j = m \\ 0, & \text{otherwise.} \end{cases} \quad (7.6)$$

Because each state represents a context, we require the transition function to ensure that we are in the right state after observing w_m : thus, we assign zero probability to all other transitions. The start function captures the probability $p(w \mid \langle \text{START} \rangle)$, and the final state function captures the probability $p(\langle \text{STOP} \rangle \mid w)$. Thus, the bigram probability of any string is computed by the product of transition scores,

$$p_2(\mathbf{w}_{1:M}) = p(w_1 \mid \langle \text{START} \rangle) \times \left(\prod_{m=2}^M p(w_m \mid w_{m-1}) \right) \times p(\langle \text{STOP} \rangle \mid w_M) \quad (7.7)$$

$$= \pi_{w_1} \times \left(\prod_{m=2}^M \delta_{q_{w_{m-1}}, w_m, q_{w_m}} \right) \times \xi_{w_M}. \quad (7.8)$$

Can you see how to construct a trigram language model in the same way?

Interpolated n-gram language model

Knight and May (2009) show how to implement an interpolated bigram/unigram language model using a WFSA. Recall that an interpolated bigram language model

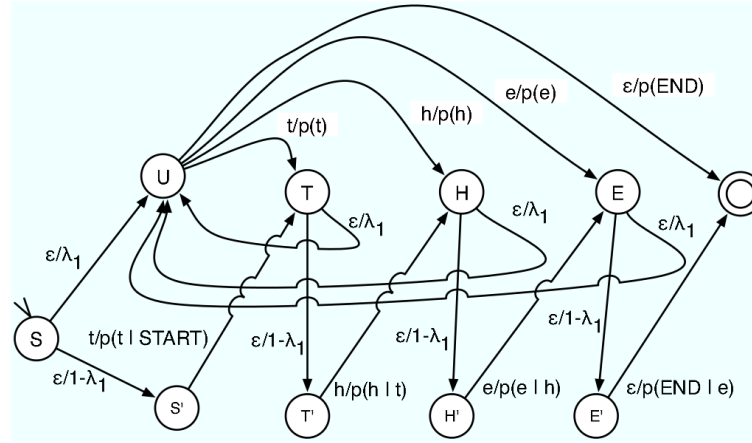


Figure 7.5: WFSA implementing an interpolated bigram/unigram language model (Knight and May, 2009). [todo: maybe redraw this for clarity?]

computes probability,

$$\hat{p}(w_m \mid w_{m-1}) = \lambda p_1(w_m) + (1 - \lambda) p_2(w_m \mid w_{m-1}), \quad (7.9)$$

with \hat{p} indicating the interpolated probability, p_2 indicating the bigram probability, and p_1 indicating the unigram probability.

Note that Equation 7.9 involves both the multiplication and addition of probabilities. Knight and May (2009) achieve this through the use of **non-determinism**. The basic idea is shown in Figure 7.5. At each of the top row of states in Figure 7.5, there are two possible ϵ -transitions, which consume no input. With score λ , we transition to the generic state U , which “forgets” the local context; transitions out of U are scored according to the unigram probability model p_1 . With score $1 - \lambda$, we transition to one of the context-remembering states, S' , T' , H' , E' . Each of these states encodes the bigram context, and outgoing transitions are scored according to the bigram probability model p_2 .

Any given path through this WFSA will have a score that multiplies together the probabilities of generating the words in the input, as well as the decisions about whether to use the unigram or bigram probability models. However, due to the non-determinism, each input string will have many possible paths to acceptance. Let’s write these paths as sequences z_1, z_2, \dots, z_M , with each $z_m \in \{1, 2\}$, indicating whether the unigram or bigram model was chosen to generating w_m .

Then the string b,a will have the following paths and scores:

$$\text{score}(1, 1, 1) = \lambda \times p_1(b) \times \lambda \times p_1(a) \times \lambda \times p_1(\langle \text{STOP} \rangle) \quad (7.10)$$

$$= \lambda^3 p_1(a) p_1(b) p_1(\langle \text{STOP} \rangle) \quad (7.11)$$

$$\text{score}(1, 1, 2) = \lambda^2(1 - \lambda) p_1(b) p_1(a) p_2(\langle \text{STOP} \rangle \mid a) \quad (7.12)$$

$$\text{score}(1, 2, 1) = \lambda^2(1 - \lambda) p_1(b) p_2(a \mid b) p_1(\langle \text{STOP} \rangle) \quad (7.13)$$

$$\text{score}(1, 2, 2) = \lambda(1 - \lambda)^2 p_1(b) p_2(a \mid b) p_2(\langle \text{STOP} \rangle \mid a) \quad (7.14)$$

$$\text{score}(2, 1, 1) = \lambda^2(1 - \lambda) p_2(b \mid \langle \text{START} \rangle) p_1(a) p_1(\langle \text{STOP} \rangle) \quad (7.15)$$

$$\text{score}(2, 1, 2) = \lambda^2(1 - \lambda) p_2(b \mid \langle \text{START} \rangle) p_1(a) p_2(\langle \text{STOP} \rangle \mid a) \quad (7.16)$$

$$\text{score}(2, 2, 1) = \lambda^2(1 - \lambda) p_2(b \mid \langle \text{START} \rangle) p_2(a \mid b) p_1(\langle \text{STOP} \rangle) \quad (7.17)$$

$$\text{score}(2, 2, 2) = (1 - \lambda)^3 p_2(b \mid \langle \text{START} \rangle) p_2(a \mid b) p_2(\langle \text{STOP} \rangle \mid a), \quad (7.18)$$

where $\langle \text{START} \rangle$ is the special start symbol and $\langle \text{STOP} \rangle$ is the special stop symbol. Each of these scores is a joint probability $p(\mathbf{w}_{1:M}, \mathbf{z}_{1:M})$; summing over them gives $\sum_{\mathbf{z}_{1:M}} p(\mathbf{w}_{1:M}, \mathbf{z}_{1:M}) = p(\mathbf{w}_{1:M})$, which is the desired marginal probability under the interpolated language model. Thus, in this case, we want not the score of the single best path, but the sum of the scores of **all** paths that accept a given input string.

7.3 Semirings

We have now seen three examples: an FSA for derivational morphology, and WFSA for edit distance and language modeling. Several things are different across these examples.

Scoring

- In the derivational morphology FSA, we wanted a boolean “score”: is the input a valid word or not?
- In the edit distance WFSA, we wanted a numerical (integer) score, with lower being better.
- In the interpolated language model, we wanted a numerical (real) score, with higher being better.

Nondeterminism

(c) Jacob Eisenstein 2014-2015. Work in progress.

- In the derivational morphology FSA, we accept if there is any path to a terminating state.
- In the edit distance WFSA, we want the score of the single best path.
- In the interpolated language model, we want to sum over non-deterministic choices.

Semiring notation allows us to combine all of these different possibilities into a single formalism.

Formal definition

A semiring is a system $(\mathbb{K}, \oplus, \otimes, \bar{0}, \bar{1})$

- \mathbb{K} is the set of possible values, e.g. $\{\mathbb{R}_+ \cup \infty\}$, the non-negative reals union with infinity
- \oplus is an addition operator
- \otimes is a multiplication operator
- $\bar{0}$ is the additive identity
- $\bar{1}$ is the multiplicative identity

A semiring must meet the following requirements:

- $(a \oplus b) \oplus c = a \oplus (b \oplus c)$, $(\bar{0} \oplus a) = a$, $a \oplus b = b \oplus a$
- $(a \otimes b) \otimes c = a \otimes (b \otimes c)$, $a \otimes \bar{1} = \bar{1} \otimes a = a$
- $a \otimes (b \oplus c) = (a \otimes b) \oplus (a \otimes c)$, $(a \oplus b) \otimes c = (a \otimes c) \oplus (b \otimes c)$
- $a \otimes \bar{0} = \bar{0} \otimes a = \bar{0}$

Semirings of interest :

where $\oplus_{\log}(a, b)$ is defined as $\log(e^a + e^b)$.

Semirings allow us to compute a more general notion of the “shortest path” for a WFSA.

- Our initial score is $\bar{1}$
- When we take a step, we use \otimes to combine the score for the step with the running total.
- When nondeterminism lets us take multiple possible steps, we combine their scores using \oplus .

Name	\mathbb{K}	\oplus	\otimes	$\bar{0}$	$\bar{1}$	Applications
Boolean	$\{0, 1\}$	\vee	\wedge	0	1	identical to an unweighted FSA
Probability	\mathbb{R}_+	+	\times	0	1	sum of probabilities of all paths
Log-probability	$\mathbb{R} \cup -\infty \cup \infty$	\oplus_{\log}	+	$-\infty$	0	log marginal probability
Tropical	$\mathbb{R} \cup -\infty \cup \infty$	min	+	∞	0	best single path

Example Let's see how this works out for our language model example.

$$\begin{aligned} \text{score}(\{a, b, a\}) = \bar{1} \otimes (\lambda \otimes \mathbf{p}_2(a|*) \oplus (1 - \lambda) \otimes \mathbf{p}_1(a)) \\ \otimes (\lambda \otimes \mathbf{p}_2(b|a) \oplus (1 - \lambda) \otimes \mathbf{p}_1(b)) \\ \otimes (\lambda \otimes \mathbf{p}_2(a|b) \oplus (1 - \lambda) \otimes \mathbf{p}_1(a)) \end{aligned}$$

Now if we plug in the **probability semiring**, we get

$$\begin{aligned} \text{score}(\{a, b, a\}) = 1 \times (\lambda \mathbf{p}_2(a|*) + (1 - \lambda) \mathbf{p}_1(a)) \\ \times (\lambda \mathbf{p}_2(b|a) + (1 - \lambda) \mathbf{p}_1(b)) \\ \times (\lambda \mathbf{p}_2(a|b) + (1 - \lambda) \mathbf{p}_1(a)) \end{aligned}$$

But if we plug in the **log probability semiring**, we need the edge weights to be equal to $\log \mathbf{p}_1$, $\log \mathbf{p}_2$, $\log \lambda$, and $\log(1 - \lambda)$. Then we get:

$$\begin{aligned} \text{score}(\{a, b, a\}) = 0 + \log (\exp(\log \lambda + \log \mathbf{p}_2(a|*)) + \exp(\log(1 - \lambda) + \log \mathbf{p}_1(a))) \\ + \log (\exp(\log \lambda + \log \mathbf{p}_2(b|a)) + \exp(\log(1 - \lambda) + \log \mathbf{p}_1(b))) \\ + \log (\exp(\log \lambda + \log \mathbf{p}_2(a|b)) + \exp(\log(1 - \lambda) + \log \mathbf{p}_1(a))) \\ = 0 + \log(\lambda \mathbf{p}_2(a|*) + (1 - \lambda) \mathbf{p}_1(a)) \\ + \log(\lambda \mathbf{p}_2(b|a) + (1 - \lambda) \mathbf{p}_1(b)) \\ + \log(\lambda \mathbf{p}_2(a|b) + (1 - \lambda) \mathbf{p}_1(a)), \end{aligned}$$

which is exactly equal to the log of the score from the probability semiring.

- The score on any specific path will be the semiring **product** of all steps along the path.

- The score of any input will be the semiring **sum** of the scores of all paths that successfully process the input.
- What happens if we use the tropical semiring?

7.4 Finite state transducers

Finite state acceptors can determine whether a string is in a language, and weighted finite state acceptors can compute a score for every string from a given alphabet. We now consider a family of automata which can **transduce** one string into another. Formally, finite state transducers (FSTs) define **regular relations** over pairs of strings. We can think of them in two different ways:

- **Recognizer:** An FST accepts a pair of strings (input and output) if the pair is in the regular relation defined by the transducer.
- **Translator:** An FST takes an input string, and returns an output, such that the input/output pair is in the regular relation.

Like FSAs, finite-state transducers are defined as tuples. In this case, we define $M = \langle Q, \Sigma, \Delta, q_0, F, \delta, \sigma \rangle$, including:

- a finite set of states $Q = \{q_0, q_1, \dots, q_n\}$;
- the finite alphabets Σ for input symbols and Δ for output symbols;
- an initial state $q_0 \in Q$, and a set of final states $F \subseteq Q$;
- a transition function $\delta : \langle Q \times \Sigma^* \rangle \rightarrow \langle Q \times \Delta^* \rangle$.

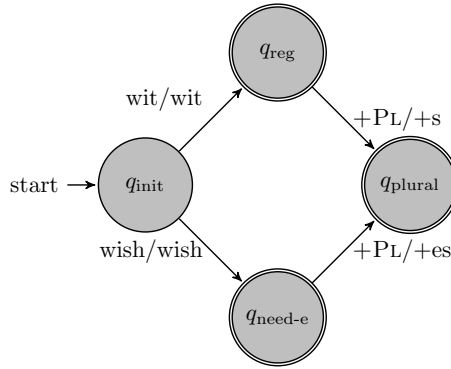


Figure 7.6: A finite state transducer for pluralizing English words.

Example Consider the following FST, shown in Figure 7.6, which performs **pluralization** of some English words:

$$Q = \{q_0, q_{\text{regular}}, q_{\text{needs-e}}, q_{\text{pluralized}}\} \quad (7.19)$$

$$N = \{aardvark, \dots, wish, wit, \dots, zyzzzyva^2\} \text{ (the set of all English nouns)} \quad (7.20)$$

$$\Sigma = N \cup \{+PL\} \quad (7.21)$$

$$\Delta = N \cup \{+s, +es\} \quad (7.22)$$

$$q_0 = q_0 \quad (7.23)$$

$$F = \{q_{\text{regular}}, q_{\text{needs-e}}, q_{\text{pluralized}}\} \quad (7.24)$$

$$\begin{aligned} \delta = \{ & (\langle q_0, aardvark \rangle \rightarrow \langle q_{\text{regular}}, aardvark \rangle), \\ & (\langle q_0, wish \rangle \rightarrow \langle q_{\text{needs-e}}, wish \rangle), \\ & (\langle q_0, wit \rangle \rightarrow \langle q_{\text{regular}}, wit \rangle), \\ & \dots \\ & (\langle q_{\text{regular}}, +PL \rangle \rightarrow \langle q_{\text{pluralized}}, +s \rangle) \\ & (\langle q_{\text{needs-e}}, +PL \rangle \rightarrow \langle q_{\text{pluralized}}, +es \rangle) \end{aligned} \quad (7.25)$$

This machine will accept the pairs $\langle wit+PL, wits \rangle$, $\langle wish+PL, wishes \rangle$, $\langle wit, wit \rangle$, but not the pairs $\langle wit+PL, wites \rangle$, $\langle wish+PL, wishs \rangle$, $\langle wish+PL, wish \rangle$. Thus, it correctly handles a small part of English orthography for pluralization; with a different word list, it could also be used to conjugate verbs to third-person singular. Consider how you might modify this FST to perform lemmatization.

Non-determinism Unlike non-deterministic finite state acceptors, not all non-deterministic finite state transducers (NFSTs) can be determinized. However, spe-

cial subsets of NFSTs called **subsequential** transducers can be determinized efficiently (see 3.4.1 in Jurafsky and Martin (2009)).

7.5 Weighted FSTs

Weights can be added to FSTs in much the same way as they are added to FSAs. For any pair $\langle q \in Q, s \in \Sigma^* \rangle$, we have a set of possible transitions, $\langle q \in Q, t \in \Delta^*, \omega \in \mathbb{K} \rangle$, with a weight ω in the domain defined by the semiring. Table 7.1 shows the relationship between FSAs, FSTs, and their weighted generalizations.

	acceptor	transducer
unweighted	FSA: $\Sigma^* \rightarrow \{0, 1\}$	FST: $\Sigma^* \rightarrow \Sigma^*$
weighted	WFSA: $\Sigma^* \rightarrow \mathbb{K}$	WFST: $\Sigma^* \rightarrow \langle \Sigma^*, \mathbb{K} \rangle$

Table 7.1: A unified view of finite state automata

Example In section 7.2, we saw how to build an FSA that would compute the edit distance from any single word. With WFSTs, we can build a general edit distance computer, which computes the edit distance between any **pair** of words.

- $Q_0 \xrightarrow[a]{a} Q_0 : 0$
- $Q_0 \xrightarrow[\epsilon]{a} Q_0 : 1$
- $Q_0 \xrightarrow[a]{\epsilon} Q_0 : 1$

The shortest path for a pair of strings $\langle s, t \rangle$ in this transducer has a score equal to the minimum edit distance between the strings (in the tropical semiring). We can think of each path as defining a potential **alignment** between s and t . That is, there are many ways to transduce *she* into *he*; in the minimum edit distance path, we have the alignment $\langle s, \epsilon \rangle, \langle h, h \rangle, \langle e, e \rangle$.

Operations on FSTs

FSTs are:

- Closed under **union**. If T_1 recognizes the relation R_1 and T_2 recognizes the relation R_2 , then there exists an FST that recognizes the relation $R_1 \cup R_2$.
- Closed under **inversion**. If T_1 recognizes the relation $R_1 = \{s_i, t_i\}_i$, then there exists an FST that recognizes the relation defined by $\{t_i, s_i\}_i$, effectively switching the inputs and outputs.
- Closed under **projection**. If T_1 recognizes the relation $R_1 = \{s_i, t_i\}_i$, then there exist FSTs that recognize the relations defined by $\{s_i, \epsilon\}_i$ and $\{\epsilon, t_i\}_i$. Note that these relations ignore either the input or the output, and so are equivalent to finite state acceptors (FSAs).
- Not closed under **difference**, **complementation**, and **intersection**;
- Closed under **composition**, as described below.

FST composition is the basis for implementing the noisy channel model in FSTs, and can be used to support dozens of cool applications. Through composition, we can create finite state **cascades** that link together several simple models; closure guarantees that the resulting model is still a WFST.

Finite state composition

Suppose we have a transducer T_1 from Σ^* to Γ^* , and another transducer T_2 from Γ^* to Δ^* . Then the composition $T_1 \circ T_2$ is an FST from Σ^* to Δ^* . More formally,

Unweighted definition iff $\langle x, z \rangle \in T_1$ and $\langle z, y \rangle \in T_2$, then $\langle x, y \rangle \in T_1 \circ T_2$.

Weighted definition

$$(T_1 \circ T_2)(x, y) = \bigoplus_{z \in \Sigma^*} T_1(x, z) \otimes T_2(z, y) \quad (7.26)$$

Note that weighted composition in the Boolean semiring is identical to unweighted composition.

Designing algorithms for automatic FST composition is relatively straightforward if there are no epsilon transitions; otherwise it's more challenging (Allauzen et al., 2009). Luckily, software toolkits like OpenFST take care of this for you.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Example

- $T_1 : Q_0 \xrightarrow[a]{x} Q_0, Q_0 \xrightarrow[b]{y} Q_0$
- $T_2 : Q_1 \xrightarrow{a} Q_1, Q_1 \xrightarrow{b} Q_2, Q_2 \xrightarrow{b} Q_2$
- $T_1 \circ T_2 : Q_1 \xrightarrow{x} Q_1, Q_1 \xrightarrow{y} Q_2, Q_2 \xrightarrow{y} Q_2$

For simplicity T_2 is written as a finite-state acceptor, not a transducer. Acceptors are a special case of transducers, where the output alphabet is $\Delta = \{\epsilon\}$.

7.6 Applications of finite state composition

Edit distance

Consider the general edit distance computer developed in section 7.5. It assigns scores to pairs of strings. If we compose it with an FSA for a given string (e.g., *tech*), we get a WFSA, who assigns score equal to the minimum edit distance from *tech* for the input string.

- Composing an FST with a FSA yields a FSA.
- A very useful design pattern is to build a **decoding** WFSA by composing a general-purpose WFST with an unweighted FSA representing the input.
- The best path through the resulting WFSA will be the minimum cost / maximum likelihood decoding.

Transliteration

English is written in a Roman script, but many languages are not. **Transliteration** is the problem of converting strings between scripts. It is especially important for names, which don't have agreed-upon translations.

A simple transliteration system can be implemented through the noisy-channel model.

- T_1 is an English character model, implemented as a transducer so that strings are scored as $\log p_r(c_1, c_2, \dots, c_M)$.
- T_2 is a character-to-character transliteration model. This can be based on explicit rules,³ or on conditional probabilities $\log p_t(c^{(f)} | c^{(r)})$.

³http://en.wikipedia.org/wiki/Romanization_of_Russian

- T_3 is an acceptor for a given string that is to be transliterated.

The machine $T_1 \circ T_2 \circ T_3$ scores English character strings based on their orthographic fluency (T_1) and adequacy (T_2).

Suppose you were given an Roman-script character model and a set of foreign-script strings, but no equivalent Roman-script strings. How would you use EM to learn a transliteration model?

Knight and May (2009) provide a more complex transliteration model, which transliterates between Roman and Katakana scripts, using a deep cascade that includes models of the underlying phonology. In their model,

Word-based translation

Machine translation can be implemented as a finite-state cascade. A simple approach is to compose three automata:

- T_1 is a language model, implemented as a transducer, where every path inputs and outputs the same string, with a score equal to $\log p(w_1, w_2, \dots, w_M)$. This model's responsibility is to tell us that $p(\text{Coffee black me pleases much}) \ll p(\text{I like black coffee a lot})$.
- T_2 is the translation machine. It contains a single state, and every transition takes a word from the source language and outputs a word in the target language. The weights are typically set to $p(w^{(t)} \mid w^{(s)})$. This model should assign a high probability to $p(\text{cafe} \mid \text{coffee})$, and a low probability to $p(\text{cafe} \mid \text{tea})$.

Suppose we are translating Spanish to English. Then T_1 maps from English to English, since it is a language model in English; T_2 maps from English to Spanish. By the definition of finite state composition (Equation 7.26), the scores of the paths through these two transducers will be combined with the \otimes operator; in the probability semiring, this means we will compute $p(w^{(e)})p(w^{(s)} \mid w^{(e)}) = p(w^{(s)}, w^{(e)})$.

- T_3 is a deterministic finite-state acceptor, which accepts only the sentence to be translated. By composing $T_1 \circ T_2 \circ T_3$, we get a weighted finite-state acceptor for sentences in the target language (in our example, English).

Recall that the composition $T_1 \circ T_2$ represents the joint probability $p(w^{(s)}, w^{(e)})$. The effect of T_3 is to “lock” $w^{(s)}$ to the sentence to be translated. The shortest

path in the composed machine $T_1 \circ T_2 \circ T_3$ thus computes,

$$\hat{w}^{(e)} = \arg \max_{w^{(e)}} p(w^{(s)}, w^{(e)}) \quad (7.27)$$

$$= \arg \max_{w^{(e)}} p(w^{(e)} \mid w^{(s)}), \quad (7.28)$$

which is the maximum-likelihood translation.

- Finally, note that we will need to allow ϵ -transitions in the translation model to handle cases like the translation of *mucho* to *a lot*. This introduces non-determinism to the finite-state cascade; again, we can think of this in terms of possible **alignments** between the source and target languages. The shortest-path algorithm computes the maximum likelihood translation while implicitly summing over all alignments.

7.7 Discriminative structure prediction

Now suppose we would like to use perceptron to learn to perform morphological segmentation. Imagine we are given a set of words $x_{1:N}$ and their true segmentations $y_{1:N}$. We would like to use perceptron to learn the weights of a WFST. How can we do it?

Recall that perceptron relies on computing a feature function $f(x, y)$. We will make this feature vector exactly equal to the finite-state transitions taken in the shortest-path transduction of x to y . That is, each potential transition $(Q_i, \omega) \rightarrow Q_o$ corresponds to some entry j in the vector $f(x, y)$, and the value $f_j(x, y)$ is equal to the number of times that transition was taken. Although FSTs can manipulate arbitrarily long strings, there will still be only a finite number of possible transitions, since both the state space and the alphabet are finite. The scores for these transitions can then be formed into the vector of weights θ , so that the score of the best path from x to y can be represented as the inner product $\theta^\top f(x, y)$.

Let these transitions be represented in the weighted FST T . Given an instance x , we build a chain acceptor A_x . By composing T and A_x , we obtain a WFSA in which the shortest path corresponds to the prediction \hat{y} , and the transitions on this path are the feature vector $f(x, \hat{y})$. We then compute the score of the best scoring path for accepting the true y segmentation in this machine; the transitions on this path form the feature vector $f(x, y)$. Given these two feature vectors, the perceptron update is as usual: $\theta^{(t+1)} \leftarrow \theta^{(t)} + f(x, y) - f(x, \hat{y})$. Weight averaging and passive-aggressive can be applied here, just as they were applicable in straightforward classification.

But unlike classification, we have now learned a function for making predictions over an **infinite set of labels**: all possible morphological segmentations for all possible words. We were able to do this by designing a feature function that shares features across different labels: if y and \hat{y} are nearly the same, then they will involve many of the same finite-state transitions, and so the feature vector $f(x, y)$ and $f(x, \hat{y})$ will be nearly the same too. This is a powerful idea that will enable us to apply the tools of classification to a huge range of problems in language technology, including part-of-speech tagging, parsing, and even machine translation.

Chapter 8

Part-of-speech tagging

Words can be grouped into rough classes based on syntax.

- Why is *colorless green ideas sleep furiously* more acceptable than *ideas colorless furiously green sleep*?
- Why is *teacher strikes idle children* ambiguous?

In both examples, word classes can provide an explanation.

- Word classes have strong ordering constraints:
 - J J N V R is relatively likely. This is the tag sequence for *colorless green ideas sleep furiously*. The abbreviation *J* means adjective, *N* means noun, *V* means verb, and *R* means adverb.
 - N J R J V is very unlikely in English. Do you see why?
- Ambiguity about word class leads to very different interpretations:
 - (8.1) *teacher/N strikes/N idle/V children/N*
 - (8.2) *teacher/N strikes/V idle/J children/N* (ouch!)

So clearly we have intuitions about a few parts-of-speech already: noun, verb, adjective, adverb. Jurafsky and Martin (2009) describe these as the four major **open** word classes, although apparently not all languages have all of them.

What other parts of speech are there?

- The Penn Treebank defines a set of 45 POS tags for English.¹

¹<http://www.comp.leeds.ac.uk/ccalas/tagsets/upenn.html>

- The Brown corpus defines a set of 87 POS tags for English.²
- Petrov et al. (2012) define a “universal” set of 12 tags, which are supposed to apply across many languages.

To understand the linguistic differences between these tagsets, let’s look at an example:

(8.3) My name is Ozymandias, king of kings:
Look on my works, ye Mighty, and despair!

The part-of-speech tags for this couplet from Ozymandias are shown in Table 8.1.

Tagset granularity

All tagsets distinguish basic categories like nouns, pronouns, verbs, adjectives, and punctuation. The Brown tagset includes a number of fine-grained distinctions:

- specific tags for the *be*, *do*, and *have* verbs, which the other two tagsets just lump in with other verbs;
- distinct tags for possessive determiners (*my name*) and possessive pronouns (*mine*);
- distinct tags for the third-person singular pronouns (e.g., *it*, *he*) and other pronouns (e.g., *they*, *we*, *I*).

In contrast, the Universal tagset aggressively groups categories that are distinguished in the other tagsets:

- all nouns are grouped, ignoring number and the proper/common distinction (see below);
- all verbs are grouped, ignoring inflection;
- preposition and postpositions are grouped as “adpositions”;
- all punctuation is grouped;
- coordinating and subordinating conjunctions (e.g. *and* versus *that*) are grouped.

²<http://www.comp.leeds.ac.uk/ccalas/tagsets/brown.html>

	Brown	PTB	Universal
My	possessive determiner (DD\$)	possessive pronoun (PRP\$)	pronoun (PRON)
name	noun, singular, common (NN)	NN	NOUN
is	verb “to be” 3rd person, singular (BEZ)	verb 3rd person, singular (VBZ)	VERB
Ozymandias	proper noun, singular (NP)	proper noun, singular (NNP)	NOUN
,	comma (,)	comma (,)	punctuation (.)
king	NN	NN	NOUN
of	preposition (IN)	preposition (IN)	adposition (ADP)
kings	noun, plural, common (NNS)	NNS	NOUN
:	colon (:)	mid-sentence punc (:)	.
Look	verb, base: uninflected present, imperative, or infinite (VB)	VB	VERB
on	IN	IN	ADP
my	DD\$	PRP\$	PRON
works	NNS	NNS	NOUN
ye	personal pronoun, nominative, non 3S (PPSS)	personal pronoun, nominative (PRP)	PRON
mighty	adjective (JJ)	JJ	adjective (ADJ)
,	comma (,)	comma (,)	punctuation (.)
and	coordinating conjunction (CC)	CC	conjunction (CONJ)
despair	VB	VB	VERB

Table 8.1: Part-of-speech annotations from three tagsets for the first couple of the poem Ozymandias.

The Penn Treebank strikes a middle ground between these two relative extremes. But which is right? It depends. The Brown tags can be useful for certain applications, and they may have strong tag-to-tag relations that make tagging easier, as described in the next chapter). But they are more expensive to annotate. The Universal tags are intended to generalize across many languages and many types of text, and should be easier to annotate.

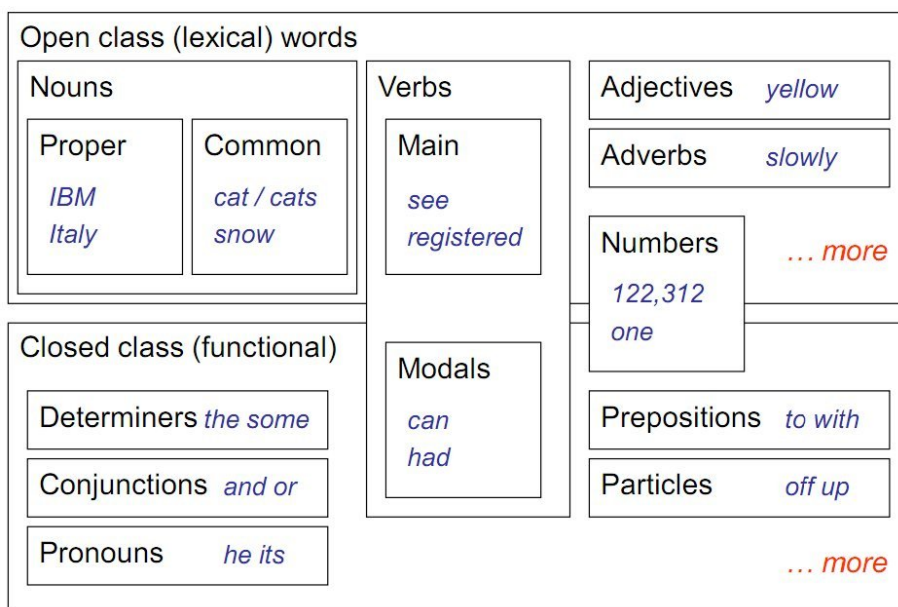


Figure 8.1: [todo: attribution?]

8.1 Details about parts-of-speech

As usual, Bender (2013) provides a useful linguistic perspective.

- **Nouns** describe entities and concepts
 - **Proper nouns** name specific people and entities: *Georgia Tech*, *Janet*, *Buddhism*. In English, proper nouns are usually capitalized. The Penn Treebank (PTB) tags are: NNP (singular), NNPS (plural).
 - **Common nouns** cover all other nouns. In English, they are often preceded by determiners, e.g. *the book*, *a university*, *some people*. Common nouns decompose into two main types:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- * **Count nouns** have a plural and need an article in the singular, *dogs, the dog*;
- * **Mass nouns** don't have a plural and don't need an article in the singular:

(8.4) *snow is cold*

(8.5) *gas is expensive*

– **Pronouns** refer to specific noun phrases or entities or events.

- * **Personal pronouns** refer to people or entities: *you, she, I, it, me*. The PTB tag is PRP.
- * **Possessive pronouns** are pronouns that indicate possession: *your, her, my, its, one's, our*. The PTB tag is PRP\$.
- * **Wh-pronouns** (WP) are used in question forms, and as relative pronouns:

(8.6) *Where are you going?*

(8.7) *The girl who played with fire.*

Unlike other nouns, the set of possible pronouns cannot be expanded. It is a **closed class**. Can you think of other closed class word groups?

- **Verbs** describe activities, processes, and events. For example, *eat, write, sleep* are all verbs.
 - The Penn Treebank differentiates verbs by morphology: VB (infinitive), VBD (past), VBG (present participle), VBN (past participle), VBZ (present 3rd person singular), VBP (present, non-3rd person singular).
 - **Modals** are a closed subclasses of verbs, such as (*should, can, will, must*). They get PTB tag MD.
 - The verb *to be* requires special treatment, as it must appear with a predicative adjective or noun, e.g.

(8.8) *She is hungry.*

(8.9) *We are Georgians.*

The verbs *is* and *are* in these cases are called **copula**. The Brown Tagset distinguishes copula, but the PTB does not. More generally, in **light verb** constructions, the meaning is largely shaped by a predicative adjective, e.g. *he got fired*, [**todo: more examples**].

- **Auxiliary** verbs include *be, have, will*, which form complex tenses in English, e.g. *we will have done it twice*. Recall from chapter 6 that English makes extensive use of auxiliary verbs to determine the tense, while other languages, such as French, rely more on morphology.
 - * Another auxiliary verb is *do*, as used in questions and negation, e.g.
 - (8.10) *Did you eat yet?*
 - (8.11) *We did not take your bagels.*
 - * The Brown corpus has special tags for HAVE and DO, but the PTB does not.
- **Adjectives** describe properties of entities: in the Ozymandias example, the adjectives include *antique, vast, trunkless*. In English, adjectives can be used in two ways:
 - **Attributive**: *an antique land*;
 - **Predicative**: *the land was antique*.

Adjectives may be **gradable**, meaning that they have a **comparative form** (e.g., *bigger, smellier*) **superlative form** (*biggest, smelliest*). Adjectives like *antique* are not gradable.

 - With *big*, we can move to comparative form by adding the suffix *-est*. This is an example of agglutinative morphology, since the comparative morpheme is added to the stem as an affix. But there are adjectives in English where the relationship between the base and comparative forms is not agglutinative, but fusional. One example *good, better, best*; can you think of any others?
 - The PTB distinguishes these forms with three tags: JJ, JJR, JJS.
- **Adverbs** describe properties of events.
 - **Manner**: *slowly, slower, fast, hesitantly*
 - **Degree**: *extremely, very, highly*
 - Adverbs may be directional or locative. In the following examples, the bolded words are all adverbs.

(8.12) *She lives **downstairs**.*

(8.13) *I study **here**.*

(8.14) Go *left* at the first traffic light.

- Adjectives also include temporal information, such as *yesterday*, *Monday*, and *soon*.
- Besides verbs, adverbs may also modify sentences, adjectives, or other adverbs.

(8.15) *Apparently*, the *very* ill man walks *extremely slowly*.

In this example, *very* modifies the adjective *ill*, *slowly* modifies the verb *walks*, *extremely* modifies the adverb *slowly*, and *apparently* modifies the entire sentence that follows it.

- Like adjectives, adverbs may also be gradable. The PTB distinguishes graded adjectives with the tags RB, RBR, RBS.
- **Prepositions** are a closed class of words that can come before noun phrases, forming a prepositional phrase that relates the noun phrase to something else in the sentence.
 - I eat sushi *with* soy sauce. The prepositional phrase **attaches** to the noun *sushi*.
 - I eat sushi *with* chopsticks. The prepositional phrase here attaches to the verb *eat*.

The preposition *To* gets its own tag TO, because it forms the **infinitive** with bare form verbs (VB), e.g. *I want to eat*. All other prepositions are tagged IN in the PTB.

- **Coordinating conjunctions** (PTB tag: CC) join two elements,

(8.16) vast *and* trunkless legs

(8.17) She plays backgammon *or* she does homework.

(8.18) She eats *and* drinks quickly.

(8.19) Sandeep lives north of Midtown *and* south of Buckhead.

(8.20) Max cooked, *and* Abigail ate, all the pizza.

- **Subordinating conjunctions** introduce a subordinate clause, e.g.

(8.21) She thinks *that* Chomsky is wrong about language models.

The PTB tag here is IN.

- **Particles** are words that come with verbs and can change their meaning to a new **phrasal verb**, e.g.,

(8.22) *Come **on**.*

(8.23) *He brushed himself **off***

(8.24) *Let's check **out** that new restaurant.*

Particles are a closed class, and are tagged RP in the PTB.

- **Determiners** (PTB tag: DT) are a closed class of words that precede noun phrases.
 - Articles: *the, an, a*
 - Demonstratives: *this, these, that*
 - Quantifiers: *some, every, few*
 - Wh-determiners: e.g., ***Which** bagel should I choose?, Do you know **when** it will be ready?*
- **Oddballs**
 - **Existential there**, e.g. *There is no way out of here*, gets its own tag, EX.
 - So does the possessive ending 's, which is POS. Recall that possessive pronouns don't have this ending, so they get a special tag, PRP\$.
 - Other special tags are reserved for numbers (CD), list items (LS), commas (,), and other non-alphabetic symbols.

8.2 Part of speech tagging

Part of speech tags relate to many other linguistic phenomena:

- Lexical semantics: *can/V* vs *can/N*, *teacher strikes children*, etc
- Pronunciation: *inSULT/V* vs *INsult/N*, *conTENT/J* vs *CONtent/N*
- Translation: *park/v* → *garer*, *park/N* → *parque*
- NP chunking: `grep {JJ | NN}* {NN | NNS}`

This means that part-of-speech tagging is a useful preprocessing step for downstream applications. So the logical next question is: how can we build an automatic POS tagger?

- Observation 1: it's easy.
 - In English, 60% of word types have only one possible POS tag.
 - If you choose the majority POS tag for each token, you get 90% right.
- Observation 2: it's not easy: a few words have a lot of possible POS tags.

(8.25) *We're taking it **back**/RB.*

(8.26) *The bar is in the **back**/NN.*

(8.27) *Go **back**/RP home. [todo: adverb?]*

(8.28) *He **backs**/VBP all the conservative candidates.*

(8.29) *The **back**/JJ roads are safer.*

- Observation 3: 90% is not actually very good. $0.9^{10} \approx .3$, so you will only get 30% of ten-word sentences correct. Sentences have exponentially many possible POS sequences. For example, the four-word sentence below has 36 possible tag sequences.

VBD		VB	
VDN	VBZ	VBP	VBZ
NNP	NNS	NN	NNS
<i>fed</i>	<i>raises</i>	<i>interest</i>	<i>rates</i>

To get an idea of how we can solve part-of-speech tagging, let's look at a tougher poem, Jabberwocky:

(8.30) 'Twas brillig, and the slithy toves
 Did gyre and gimble in the wabe:
 All mimsy were the borogoves,
 And the mome raths outgrabe.

Forget *twas*. What about *slithy* and *toves*? Can you guess the part of speech? You probably don't know what these words mean, for the very good reason that they are not real words. But you might still have a good guess about their syntactic class. What information are you using to make these guesses?

- **Word identity:** you do know that *and* is CC and *the* is DET.
- **Context**

- JJ NN is a frequently observed pattern in English; So are DET JJ and DET NN.
- DET VB is rarely observed in English.

- **Morphology**

- The suffix *-s* usually indicates a noun or a verb.
- The suffix *-able* indicates an adjective — 98% of the time!
- The suffix *-ly* often indicates an adverb.
- The prefix *un-* often indicates an adjective or a verb.

But these not rules, just hints: exceptions include *uncle*, *rely*, and *stable*. We therefore need to combine these intuitions with other features of the sentence.

Let's put morphology on hold for a minute. Suppose we have an annotated corpus, with tagged sentences, $\{(\mathbf{w}_{1:N_i}, \mathbf{y}_{1:N_i})\}_{1:T}$.

- We can estimate the likelihood of a word given a tag, for example by using relative frequency estimation:

$$p(w | y) = \frac{\text{count}(w, y)}{\text{count}(y)}. \quad (8.1)$$

As in language modeling and Naïve Bayes, smoothing is usually advisable.

- Given this same annotated corpus, we can also compute $p(y_m | y_{m-1})$, which is a sort of language model over tags.

$$p(y_m | y_{m-1}) = \frac{\text{count}(y_{m-1}, y_m)}{\text{count}(y_{m-1})} \quad (8.2)$$

Let's combine these ideas via a **generative story**

- For word m , draw tag $y_m \sim \text{Categorical}(\theta_{y_{m-1}})$
- Then draw word $w_m \sim \text{Categorical}(\phi_{y_m})$

We've built a generative model that explains our observations \mathbf{w} through a bigram generative model over the tags. Under this model, we can compute,

$$p(\mathbf{y} | \mathbf{w}) \propto p(\mathbf{w}, \mathbf{y}) \quad (8.3)$$

$$= p(\mathbf{w} | \mathbf{y}) p(\mathbf{y}) \quad (8.4)$$

$$= \prod_m^M p(w_m | y_m) p(y_m | y_{m-1}) \quad (8.5)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

This is a **hidden Markov model**.

- It's **Markov** because the probability of y_m depends only on y_{m-1} and not any of the previous history.
- It's **hidden** because $\mathbf{y}_{1:M}$ is unknown when we decode a string $\mathbf{w}_{1:M}$.

Hidden Markov models are an extremely well-known concept in natural language processing. But in fact, they are just a special case of finite state transduction. Can you see how they relate?

Chapter 9

Sequence labeling

In sequence labeling, we want to assign tags to words, or more generally, to discrete elements in a sequence. There are many applications of sequence labeling in natural language processing:

- Part-of-speech tagging: *Go/V to/P Georgia/N Tech/N next/J year/N ./.*
- Named entity recognition:¹ *Go/O to/O Georgia/B-ORG Tech/I-ORG next/B-DATE year/I-DATE ./O*
- Phrase chunking: *Go/B-VP to/B-PP Georgia/B-NP Tech/I-NP next/B-NP year/I-NP ./O*

In classification, we would choose each tag independently, $y_m \perp y_n \mid w_m, \forall m, n$. But in sequence labeling, we choose the sequence of tags **jointly**. Probabilistically, we might try to choose $\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}^M} p(\mathbf{y} \mid \mathbf{w})$. As we will see later, we can also write this in the form of a linear predictor:

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}^M} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) \quad (9.1)$$

In either case, we have an immediate problem: finding the best scoring tag sequence in the set \mathcal{Y}^M . As the notation suggests, the number of possible tag sequences is exponential in the length of the sequence; we saw this in the previous chapter, where the short example *Fed raises interest rates* has 36 possible part-of-speech tag sequences! This exponential growth means we will need clever algorithms to compute $\arg \max_{\mathbf{y} \in \mathcal{Y}^M}$; we cannot possibly enumerate all possibilities.

¹These examples show **BIO** notation, in which spans such as ORG (organization) or NP (noun phrase) are delimited using prefixes B- and I-. The prefixes indicate whether each token is at the beginning or inside of the span; the tag O is reserved for tokens that are outside any span. For now, we will just think of B-ORG, I-ORG, O, etc, as separate tags; see chapter 18 for more.

9.1 Hidden Markov Models

Let's first think about tagging as a probabilistic model. Specifically, we want to maximize $p(\mathbf{y}|\mathbf{w}) \propto p(\mathbf{y}, \mathbf{w})$, where \mathbf{w} are words and \mathbf{y} are tags. This is equivalent to Naive Bayes, but for sequence labeling.

As in Naive Bayes, we define the probability distribution $p(\mathbf{w}, \mathbf{y})$ through a *generative story*,

- For word m , draw tag $y_m \sim \text{Categorical}(\lambda_{y_{m-1}})$
- Then draw word $w_m \sim \text{Categorical}(\phi_{y_m})$

Under this model, we can compute

$$p(\mathbf{y} | \mathbf{w}) \propto p(\mathbf{w}, \mathbf{y}) \quad (9.2)$$

$$= p_e(\mathbf{w} | \mathbf{y}; \phi) p_t(\mathbf{y}; \lambda) \quad (9.3)$$

$$= \prod_m^M p_e(w_m | y_m; \phi) p_t(y_m | y_{m-1}; \lambda) \quad (9.4)$$

This is a **hidden Markov model**. It's Markov because the probability of y_m depends only on y_{m-1} and not any of the previous history. It's hidden because y_m is unknown when we decode.

- The probability $p_e(w_m | y_m; \phi)$ is the **emission probability**, since the words are treated as emissions from the tags.
- The probability $p_t(y_m | y_{m-1}; \lambda)$ is the **transition probability**, since it assigns probability to each possible tag-to-tag transition.

This generative story is often represented as a graphical model. Note that although graphical models and finite-state models both use circles and arrows, the meaning is completely different; here the nodes represent random variables, and the edges represent probabilistic dependencies.

The generative story assumes that the words are conditionally independent given the tags,

$$w_n \perp \{\mathbf{w}_{m \neq n}\} \mid \mathbf{y}_n.$$

Conditional independence is not the same as independence. We do **not** have $p(w_n, w_m) = p(w_n)p(w_m)$, because the tags are related to each other. For example,

(c) Jacob Eisenstein 2014-2015. Work in progress.

suppose that (a) nouns always follow determiners, (b) *the* is always a determiner and (c) *bike* is always a noun. Then

$$Pr(W_m = the, W_{m+1} = bike) = \sum_{y_{m+1}, y_m} Pr(W_m = the, W_{m+1} = bike, y_{m+1}, y_m) \quad (9.5)$$

$$= \sum_{y_{m+1}, y_m} Pr(W_{m+1} = bike \mid y_{m+1}, y_m, W_m = the) \quad (9.6)$$

$$\times Pr(y_{m+1} \mid y_m, W_m = the) Pr(y_m \mid W_m = the) Pr(W_m = the) \quad (9.7)$$

$$= \sum_{y_{m+1}} Pr(W_{m+1} = bike \mid y_{m+1}) \quad (9.8)$$

$$\times \sum_{y_m} Pr(y_{m+1} \mid y_m) Pr(y_m \mid W_m = the) Pr(W_m = the) \quad (9.9)$$

$$= Pr(W_{m+1} = bike \mid y_{m+1} = \text{NOUN}) \times 1 \times 1 \times Pr(W_m = the) \quad (9.10)$$

$$> Pr(W_{m+1} = bike) Pr(W_m = the). \quad (9.11)$$

Since *bike* is mainly used as a noun, the conditional probability $p(bike \mid N)$ is greater than the marginal $p(bike)$.

Another way to think about independence is that if we are told one tag, it affects all of our other tagging decisions.

- For example, in the sentence *teacher strikes idle children*, we might choose tag sequence NN VBZ JJ NNS.
- But if are given $y_3 = \text{VBP}$, then suddenly $y_2 = \text{VBZ}$ looks like a bad choice because $p_T(\text{VBZ}, \text{VBP})$ is very small.
- So we might now choose $y_2 = \text{NNS}$.
- This change might cascade back to y_1 , etc (not in this case, but it could happen in theory)

A classifier-based tagger, which treated the tags as IID, might ignore these dependencies, and produce a tag sequence that contained unlikely transitions like VBZ, VBP. A better alternative might be to tag the text from left-to-right; we could then condition on the previous tag, choosing

$$y_m = \arg \max_y p_e(w_m \mid y_m) p_t(y_m \mid y_{m-1}) \quad (9.12)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

But this approach is “greedy,” and can mistakenly commit to bad tagging decisions. For example, in *teacher strikes strand children*, we might initially choose $y_2 = \text{VBZ}$, because this is more common than the noun sense of *strikes*. However, we are then stuck, because *strand* has low probability as anything but a verb, yet the verb-verb transition also has low probability. The greedy tagger is unable to recover the globally optimal sequence, NN NNS VBP NNS, without backtracking. This is why we need **joint inference** over $\mathbf{y}_{1:M}$ to find $\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} p(\mathbf{w}, \mathbf{y})$. The key challenge is to search over the exponential number of tag sequences efficiently.

9.2 Algorithms for sequence labeling

Finite state transduction

To see whether efficient joint inference is possible, we first formulate the problem in terms of finite-state transduction.

- Transducer E has one state, and transduces from tags to words. Each edge begins and ends in the same state, and has cost $\delta_{w/y, q_0 \rightarrow q_0}^{(e)} = p_e(w \mid y)$.
- Transducer T has $\#|\mathcal{T}|$ states (assuming a bigram model), and transduces tags to tags, with $\delta_{y/y, q_{y_{m-1}} \rightarrow q_{y_m}}^{(t)} = p_t(y_m \mid y_{m-1})$.

Now, recall the definition of finite state composition,

$$(T \circ E)(y, x) = \bigoplus_z T(y, z) \otimes E(z, x). \quad (9.13)$$

Since T only accepts identical tag pairs $\langle y, y \rangle$, we can ignore \bigoplus ; there is only one possible $z = y$. The result of $T \circ E$ is a WFST that transduces tags to words, with edge weights equal to,

$$\begin{aligned} \delta_{w_m/y_m, q_{y_{m-1}} \rightarrow q_{y_m}}^{(toe)} &= \delta_{w/y_m, q_0 \rightarrow q_0}^{(e)} \otimes \delta_{y_m/y_m, q_{y_{m-1}} \rightarrow q_{y_m}}^{(t)} \\ &= p(w \mid y_m) \otimes p(y_m \mid y_{m-1}) \\ &= p(w \mid y_m) p(y_m \mid y_{m-1}) \\ &= p(w, y_m \mid y_{m-1}). \end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

Suppose we wanted to work with log probabilities instead. Then,

$$\begin{aligned}
 \delta_{w_m/y_m, q_{y_{m-1}} \rightarrow q_{y_m}}^{(toe)} &= \log p(w \mid y) \\
 \delta_{y/y, q_{y_{m-1}} \rightarrow q_{y_m}}^{(t)} &= \log p(y_m \mid y_{m-1}) \\
 a \otimes b &:= a + b \\
 \delta_{w_m/y_m, q_{y_{m-1}} \rightarrow q_{y_m}} &= \log p(w_m \mid y_m) \otimes \log p(y_m \mid y_{m-1}) \\
 &= \log p(w_m \mid y_m) + \log p(y_m \mid y_{m-1}) \\
 &= \log p(w_m, y_m \mid y_{m-1}).
 \end{aligned}$$

Can you see how many states the resulting FST will have?

To **decode** an input sentence $w_{1:M}$, we compose this FST with a chain acceptor S . This FSA should accept only the sequence $w_{1:M}$. The composition $T \circ E \circ S$ yields a **trellis-shaped** weighted finite state acceptor (WFSA).

- Number of columns = M , length of input.
- Number of rows = T , number of tags.
- Edges from states $\langle m, t_1 \rangle$ to $\langle m+1, t_2 \rangle$ have the score,

$$\begin{aligned}
 \delta_{w_{m+1}/t_2, q_{m, t_1} \rightarrow q_{m+1, t_2}}^{(toe)} &= \delta_{t_2/t_2, q_{t_1} \rightarrow q_{t_2}}^{(t)} \otimes \delta_{t_2/w_{m+1}, q_0 \rightarrow q_0}^{(e)} \\
 &= Pr_T(Y_{m+1} = t_2 \mid Y_m = t_1) \quad (9.14)
 \end{aligned}$$

$$\times Pr_E(W_{m+1} = w_{m+1} \mid Y_{m+1} = t_2). \quad (9.15)$$

Each path in the trellis corresponds to a unique sequence of tags, $y_{1:M}$, and every sequence of tags has a unique path. The score of the path is equal to $p(w_{1:M}, y_{1:M})$ by construction. If we define $\oplus = \max$ (as in the tropical semiring), then the score of the semiring shortest path is equal to $\max_y p(w_{1:M}, y_{1:M})$.

The algorithmic question is: given that there are exponential number of possible paths, can we still find the best score (and therefore the best path) in polynomial time?

- **How expensive is it to construct the trellis?**
 - Generic composition is polynomial, but it depends on the vocabulary size.
 - But since we know what the trellis is supposed to look like, we can just build it directly. This requires constant time per edge, ignoring the vocabulary size.

- **How big is the trellis?** $\mathcal{O}(MT)$ states, $\mathcal{O}(MT^2)$ edges.
- **How expensive is it find the shortest path in the trellis?:**
 Generic shortest path has a time cost of $\mathcal{O}(V \log V + E)$, where V is the number of vertices and E is the number of edges. In this case, we have $V = MT$ vertices and $E = MT^2$ edges. The time cost is therefore $\mathcal{O}(MT \log MT + MT^2)$, and the space cost is $\mathcal{O}(V) = \mathcal{O}(MT^2)$.
- To summarize:
 - building the trellis is polynomial;
 - shortest path is polynomial;
 - therefore, there must be a poly-time algorithm to find the best tag sequence, despite the apparently exponential number of paths.

The Viterbi algorithm

The Viterbi algorithm is a special-purpose best-path algorithm for FSTs in the shape of a trellis. It has a time cost of $\mathcal{O}(MT^2)$ and a space cost of $\mathcal{O}(MT)$. This time cost improvement is important, because Viterbi has linear time complexity in the length of the sequence M , unlike the generic shortest-path algorithm, which is $\mathcal{O}(M \log M)$.

To understand the algorithm, note that the Markov assumption ensures that we can decompose the likelihood recursively.

$$p(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}) = p(w_M | y_M) p(y_M | y_{M-1}) p(\mathbf{w}_{1:M-1}, \mathbf{y}_{1:M-1})$$

- Given y_{m-1} , we can choose y_m without considering any other element of the history.
- Suppose we know the best path to $y_m = k$. The best path to $y_{m+1} = k'$ through $y_m = k$ must include the best path to $y_m = k$.
- Suppose we know the score (probability) of the best path to each $y_m = k$, which we write $v_m(k) = \max_{y_1 \dots y_{m-1}} p(\mathbf{w}_{1:m}, \mathbf{y}_{1:m-1}, y_m = k)$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

What is the score of the best path to $y_{m+1} = k'$?

$$v_{m+1}(k') = \max_{\mathbf{y}_{1:m}} \mathbf{p}(\mathbf{w}_{1:m+1}, \mathbf{y}_{1:m}, y_{m+1} = k') \quad (9.16)$$

$$= \mathbf{p}_e(w_{m+1} \mid y_{m+1} = k') \max_{\mathbf{y}_{1:m}} P_t(Y_{m+1} = k' \mid y_m) \mathbf{p}(\mathbf{w}_{1:m}, \mathbf{y}_{1:m}) \quad (9.17)$$

$$= \mathbf{p}_e(w_{m+1} \mid y_{m+1} = k') \max_{y_m=k} P_t(Y_{m+1} = k' \mid Y_m = k) \max_{\mathbf{y}_{1:m-1}} \mathbf{p}(\mathbf{w}_{1:m}, \mathbf{y}_{1:m-1}, y_m = k) \quad (9.18)$$

$$= \mathbf{p}_e(w_{m+1} \mid y_{m+1} = k') \max_{y_m=k} P_t(Y_{m+1} = k' \mid Y_m = k) v_m(k) \quad (9.19)$$

The base case is $v_0(\langle \text{START} \rangle) = 1$, with zero probability for everything else.

- We can generalize this recurrence using semiring notation:

$$v_{m+1}(k') = \delta_{w_{m+1}, y_{m+1}=k'}^{(e)} \otimes \left(\bigoplus_k \delta_{k \rightarrow k'}^{(t)} \otimes v_m(k) \right) \quad (9.20)$$

Then if we want to move to log-probabilities, we have

$$v_{m+1}(k') = \log \mathbf{p}_E(w_{m+1} \mid y_{m+1} = k') \otimes \left(\bigoplus_k \log \mathbf{p}_T(k \rightarrow k') \otimes v_m(k) \right) \quad (9.21)$$

$$= \log \mathbf{p}_E(w_{m+1} \mid y_{m+1} = k') + \max_k \log \mathbf{p}_T(k \rightarrow k') + v_m(k) \quad (9.22)$$

- We will frequently use a semiring in which the edge weights are log probabilities and \otimes is addition. This is partly because addition is notationally clearer than multiplication, and because in practical settings, you will use the log probabilities to avoid underflow.
- Note that we are setting $\oplus = \max$, as in the tropical semiring. This means that the score of the best tag sequence overall is $v_M(\langle \text{STOP} \rangle)$.
- To find the best tag sequence, we just need to keep back-pointers, from $v_m(k)$ to $v_{m-1}(k')$:

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$v_{m+1}(k') = \max_k \log p_E(w_{m+1} \mid Y_{m+1} = k') + \log P_T(Y_{m+1} = k' \mid Y_m = k) + v_m(k) \quad (9.23)$$

$$= \log p_E(w_{m+1} \mid y_{m+1} = k') + \left(\max_k \log P_T(Y_{m+1} = k' \mid Y_m = k) + v_m(k) \right) \quad (9.24)$$

$$b_{m+1}(k') = \arg \max_k \log p_E(w_{m+1} \mid Y_{m+1} = k') + \log P_T(Y_{m+1} = k' \mid Y_m = k) + v_m(k) \quad (9.25)$$

$$= \arg \max_k \log P_T(Y_{m+1} = k' \mid Y_m = k) + v_m(k) \quad (9.26)$$

Note that the computation of the back-pointer doesn't depend on the emission probability $p_E(w_{m+1} \mid Y_{m+1} = k')$, since Y_m is conditionally independent from w_{m+1} given Y_{m+1} .

- In the probability semiring, we had \oplus as addition; in the log-probability semiring, it was log addition. What happens if we try these addition operators? We'll see in a moment.

Example

Table 9.1: $\log p(w \mid y)$ [**todo: cannot seem to use underscores in mathmode here, some weird gb4e issue**]

<u><u>they can fish</u></u>

Table 9.2: $\log p(y_m \mid y_{m-1})$

	N	V	END
START	-1	-2	$-\infty$
N	-3	-1	-2
V	-1	-3	-2

See the slides for how the Viterbi algorithm works in this example.

The forward algorithm

In an influential survey, Rabiner (1989) defines three problems for hidden Markov models:

Decoding Find the best tags \mathbf{y} for a sequence \mathbf{w} .

Likelihood Compute the marginal probability $p(\mathbf{w}) = \sum_{\mathbf{y}} p(\mathbf{w}, \mathbf{y})$.

Learning Given only unlabeled data $\{\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_D\}$, estimate the transition and emission distributions.

The Viterbi algorithm solves the decoding problem. We'll talk about the learning problem in section 9.5. Let's now consider how to compute the likelihood $p(\mathbf{w}) = \sum_{\mathbf{y}} p(\mathbf{w}, \mathbf{y})$. Recall that the Viterbi algorithm can be written in semiring notation,

$$v_{m+1}(k') = \bigoplus_k p_E(w_{m+1} \mid Y_{m+1} = k') \otimes Pr_T(Y_{m+1} = k' \mid Y_m = k) \otimes v_m(k). \quad (9.27)$$

In the Viterbi algorithm, we used a semiring in which $a \oplus b$ is defined as $\max(a, b)$, and $a \otimes b$ is defined as $a \times b$. Now let us consider a semiring in which we redefine $a \oplus b$ to be equal to $a + b$; if the associated variables refer to probabilities then \oplus corresponds to adding variables. In this semiring, we will denote the variables on the trellis as $\alpha_m(k)$, indicating the value for tag k at token m . Let us take the inductive hypothesis that $\alpha_m(k) = p(\mathbf{w}_{1:m}, Y_m = k)$; as we will show, this enables us to recursively compute the desired joint probability, $p(\mathbf{w}_{1:M})$.

$$\alpha_{m+1}(k') = \bigoplus_k p_E(w_{m+1} \mid Y_{m+1} = k') \otimes Pr_T(Y_{m+1} = k' \mid Y_m = k) \otimes \alpha_m(k) \quad (9.28)$$

$$= \sum_k p_E(w_{m+1} \mid Y_{m+1} = k') \times Pr_T(Y_{m+1} = k' \mid Y_m = k) \times \alpha_m(k) \quad (9.29)$$

$$= \sum_k p(w_{m+1}, Y_{m+1} = k' \mid Y_m = k) \times \alpha_m(k) \quad (9.30)$$

$$= \sum_k p(w_{m+1}, Y_{m+1} = k' \mid Y_m = k) \times p(\mathbf{w}_{1:m}, Y_m = k) \quad (9.31)$$

$$= \sum_k p(\mathbf{w}_{1:m+1}, Y_{m+1} = k', Y_m = k) \quad (9.32)$$

$$= p(\mathbf{w}_{1:m+1}, Y_{m+1} = k'). \quad (9.33)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

In the base case, $\alpha_1(k) = p_E(w_1 | Y_1 = k)Pr_T(Y_1 = k | Y_0 = \langle \text{START} \rangle)$. Finally, we have,

$$\sum_k \alpha_M(k) \times p_T(\langle \text{STOP} \rangle | Y_M = k) = p(\mathbf{w}_{1:M} | Y_M = k) \times p_T(\langle \text{STOP} \rangle | Y_M = k) \quad (9.34)$$

$$= \sum_k p(\mathbf{w}_{1:M}, Y_m = k, Y_{M+1} = \langle \text{STOP} \rangle) \quad (9.35)$$

$$= p(\mathbf{w}_{1:M}). \quad (9.36)$$

This recurrence is called the **forward** algorithm. In practice, log-probabilities are more numerically stable, so we use a semiring in which,

$$a \otimes b = a + b \quad (9.37)$$

$$a \oplus b = \log(e^a + e^b). \quad (9.38)$$

This definition of semiring addition ensures that $\log p(x) \oplus \log p(y) = \log(p(x) + p(y))$.

The forward algorithm can be extended to tag trigrams in exactly the same way as the Viterbi algorithm. As in Viterbi, the time complexity is $\mathcal{O}(MK^2)$ for the tag-bigram forward algorithm, and $\mathcal{O}(MK^3)$ for the tag-trigram forward algorithm. Unlike Viterbi, there is no need to keep backpointers in the forward algorithm.

Applications of the forward algorithm

Why would we want to compute the joint probability $p(\mathbf{w}_{1:M})$? There are a few reasons:

Word class language models Remember the question of how to compute the probability of an “unseen” grammatical sentence like $p(\text{colorless green ideas sleep furiously})$. In this case, we don’t care about the specific tags, we just want to know the probability of the utterance, so we can compare it with an unseen ungrammatical sentence, $p(\text{Furiously sleep ideas green colorless})$. The forward algorithm can be used for this purpose.

Comparing HMMs Suppose we have a few HMMs, each of which could have generated the observations. If each HMM corresponds to a different explanation of the input, then we might like to know which HMM is most likely to be responsible for the observations. The forward algorithm can be used for this purpose, in a sort of sequence-level version of Naïve Bayes. This approach is sometimes used in gesture recognition (?).

Computing marginal probabilities The main practical reason for using the forward algorithm is that it can help us compute marginal probabilities for individual tags $p(y_m \mid \mathbf{w}_{1:M})$ and for tag bigrams, $p(y_m, y_{m+1} \mid \mathbf{w}_{1:M})$. These marginal probabilities are needed for learning in **conditional random fields**, described in section 9.4. Note that here we condition on the entire word sequences $\mathbf{w}_{1:M}$, so these quantities cannot be computed directly from the forward algorithm, which can only tell us $p(\mathbf{w}_{1:m} \mid y_m)$. The required probabilities are obtained by using the forward algorithm in combination with the **backward algorithm**, described in section 9.4.

9.3 Learning models of sequence labeling

In principle, we can use relative frequency estimation to compute the probabilities of emissions and transitions in the HMM.

$$\lambda_{k,k'} \triangleq \text{Pr}_T(Y_m = k' \mid Y_{m-1} = k) = \frac{\text{count}(Y_m = k', Y_{m-1} = k)}{\text{count}(Y_{m-1} = k)}$$

$$\phi_{k,i} \triangleq \text{Pr}_E(W_m = i \mid Y_m = k) = \frac{\text{count}(W_m = i, Y_m = k)}{\text{count}(Y_m = k)}$$

In practice, we need smoothing for this to work well. The same ideas from language models (chapter 5) can be applied: holding out probability mass for unseen events, and interpolating between trigrams, bigrams, and unigrams.

However, in practice, probabilistic generative models are rarely used for part-of-speech tagging or other supervised sequence labeling tasks in NLP. This is because there are two things that probabilistic generative models cannot easily give us: rich features and fine-grained context.

Rich features Recall the example of the Jabberwocky poem from chapter 8:

(9.1) 'Twas brillig, and the slithy toves
Did gyre and gimble in the wabe:
All mimsy were the borogoves,
And the mome raths outgrabe

You probably didn't know many of these words, yet it was not so hard to see what some of their tags should be. How did we do it? Recall that the HMM can incorporate two sources of information:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Word-tag probabilities, via $p_E(w_m|y_n)$.
- Local context, via $p_T(y_m|y_{m-1})$.

Local context is helpful, but the word-tag probabilities will be worthless for words like *brillig*, *slithy*, *toves*, *gyre*, etc. For these words, we might rely on guesses about the morphology. But morphological features are difficult to incorporate in a generative model, because they break the Naive Bayes assumption:

$$p(\text{mimsy}, -sy \mid JJ) \neq p(\text{mimsy} \mid JJ)p(-sy \mid JJ) \quad (9.39)$$

Similarly, in **named entity recognition**, capitalization is a particularly important feature. This is what allows us to distinguish classically ambiguous cases like *I bought an apple* and *I bought an Apple computer*.

More advanced HMMs incorporate morphological, orthographic, and typographic features by creating a more complex $p_E(w \mid y)$ emission probability. For example, the TNT Tagger took this approach, and is one of the best generative taggers (Brants, 2000). However, incorporating morphological features while preserving conditional independence is extremely challenging, making inference complex.

Fine-grained context In addition to word-internal features, we might want more fine-grained context. For example, in the PTB, *this* and *these* are both tagged DT. But *this* is likely to be followed by a singular noun NN, and *these* is likely to be followed by a plural noun NNS. So we might like to add word-context features to the probability $p(y_m \mid y_{m-1}, w_{m-1})$.

How can we incorporate these overlapping features? The solution is to build sequence labeling models based on the perceptron and logistic regression classifiers. The first model is called **structured perceptron**, since the label space consists of structures rather than individual labels (Collins, 2002). The second model is called a **conditional random field (CRF)**, due to its relation to Markov random fields (Lafferty et al., 2001). In this model, we explicitly compute $p(\mathbf{y} \mid \mathbf{w})$.

In addition to incorporating overlapping features, these models have another advantage: they are discriminative, directly maximizing the conditional probability $p(\mathbf{y} \mid \mathbf{w})$, or minimizing the perceptron loss. As in standard classification, this criterion is more closely connected to the accuracy metrics that we usually care about.

Structured perceptron

Remember the perceptron update:

$$\hat{y} = \arg \max_y \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y) \quad (9.40)$$

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} + \mathbf{f}(\mathbf{w}, y^*) - \mathbf{f}(\mathbf{w}, \hat{y}) \quad (9.41)$$

In sequence labeling, we have a **structured output** $\mathbf{y} \in \mathcal{Y}(\mathbf{w})$. Can we still apply the perceptron rule?

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) \quad (9.42)$$

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} + \mathbf{f}(\mathbf{w}, \mathbf{y}^*) - \mathbf{f}(\mathbf{w}, \hat{\mathbf{y}}) \quad (9.43)$$

This is called structured perceptron, because it learns to predict structured output \mathbf{y} . The problem is that $\arg \max_y$ must search over the entire set of structures $\mathcal{Y}(\mathbf{w})$ — and note that this set of permissible outputs depends on the input \mathbf{w} , and it is very large: $\mathcal{O}(K^M)$. What can we do?

Viterbi inference for structured perceptron

We will place a restriction on the scoring function $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y})$, which allows us to apply the Viterbi algorithm to find $\arg \max_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y})$!

- Specifically, we require $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_m \boldsymbol{\theta}^\top \mathbf{f}_m(\mathbf{w}, y_m, y_{m-1}, m)$
- That is, the global score must be a **sum of local scores**.
- The local scores can consider any part of the observation, but must only consider adjacent elements in the label.

To apply Viterbi to structured perceptron, we have

$$\begin{aligned} v_m(k) &= \bigoplus_{k'} \boldsymbol{\theta}^\top \mathbf{f}_m(\mathbf{w}, k, k', m) \otimes v_{m-1}(k') \\ &= \max_{k'} \boldsymbol{\theta}^\top \mathbf{f}_m(\mathbf{w}, k, k', m) + v_{m-1}(k') \\ b_m(k) &= \arg \max_{k'} \boldsymbol{\theta}^\top \mathbf{f}_m(\mathbf{w}, k, k', m) + v_{m-1}(k') \end{aligned}$$

Suppose we want to apply this to POS tagging? What features might we want? Here are some:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Word-tag features, e.g. $\langle W : \text{slithy}, \text{JJ} \rangle$
- Adjacent tag-tag features, e.g. $\langle T : \text{JJ}, \text{NNS} \rangle$
- Suffix-tag features, e.g., $\langle M : \text{-es}, \text{NNS} \rangle$
- Previous-word features, e.g., $\langle P_1 : \text{the}, \text{JJ} \rangle$
- Next-word features, e.g., $\langle N_1 : \text{slithy}, \text{DT} \rangle$
- Note that we can consider arbitrarily distant words, e.g. $\langle Y_m, W_{m-15} \rangle$, because this still fits in the constraint, $\theta^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_m \theta^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$.

So to compute $v_m(k)$, we have to iterate over all $y_{m-1} = k'$,

- find the features $\mathbf{f}(\mathbf{x}_{1:M}, y_m = k, y_{m-1} = k', m)$,
- compute the inner product $\theta^\top \mathbf{f}(\mathbf{x}_{1:M}, y_m = k, y_{m-1} = k', m)$,
- add it to $v_{m-1}(k')$
- take the max over all k'

This only works because of the assumption that the feature function decomposes over local parts of the sequence! If we wanted a feature that considered arbitrary parts of the tag sequence, there would be no way to incorporate it into the recurrence relation.

Example

$\mathbf{w} = \dots \text{and the slithy toves}$
 $\mathbf{y} = \dots \text{CC DT JJ NNS}$

Then we can characterize the tagging DT JJ NNS of the text the slithy toves (from Jabberwocky) in terms of the following features:

$$\begin{aligned} \mathbf{f}(\text{the slithy toves}, \text{DT JJ NNS}) = & \{ \langle W : \text{the}, \text{DT} \rangle, \langle M : \emptyset, \text{DT} \rangle, \langle T : \langle \text{START} \rangle, \text{DT} \rangle \\ & \langle W : \text{slithy}, \text{JJ} \rangle, \langle M : \text{-thy}, \text{JJ} \rangle, \langle T : \text{DT}, \text{JJ} \rangle \\ & \langle W : \text{toves}, \text{NNS} \rangle, \langle M : \text{-es}, \text{NNS} \rangle, \langle T : \text{JJ}, \text{NNS} \rangle \\ & \langle T : \text{NNS}, \langle \text{STOP} \rangle \rangle \} \end{aligned}$$

If this is the correct tagging, then we hope to learn a set of weights \mathbf{w} such that $\theta^\top \mathbf{f}(\text{the slithy toves}, \text{DT JJ NNS})$ is larger than the scores for other tag sequences, such as $\theta^\top \mathbf{f}(\text{the slithy toves}, \text{DT NN VBZ})$.

Learning

If we define $\mathbf{f}(\mathbf{w}_{1:M}, y_m, y_{m-1}, m) = \{\langle W : w_m, y_m \rangle, \langle T : y_{m-1}, y_m \rangle\}$, then the expressiveness of the model is identical to the HMM. By setting the weights of these features to the log of their maximum-likelihood estimates,

$$\begin{aligned} w_{\langle W : w_m, y_m \rangle} &= \log \text{count}(w_m, y_m) - \log \text{count}(y_m) \\ w_{\langle T : y_{m-1}, y_m \rangle} &= \log \text{count}(y_{m-1}, y_m) - \log \text{count}(y_{m-1}), \end{aligned}$$

we exactly recover the hidden Markov model.

To use overlapping features and to get the advantages of error-driven learning, we can use perceptron updates. Learning is exactly the same as the non-structured perceptron from chapter 2:

$$\hat{\mathbf{y}} \leftarrow \arg \max_{\mathbf{y}} (\boldsymbol{\theta}^{(t)})^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) \quad (9.44)$$

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} + \mathbf{f}(\mathbf{x}, \mathbf{y}) - \mathbf{f}(\mathbf{x}, \hat{\mathbf{y}}). \quad (9.45)$$

The key difference is that instead of computing $\hat{\mathbf{y}}$ by enumerating the entire set \mathcal{Y} , we use the Viterbi algorithm to search this set efficiently.

As before, weight averaging is crucial to get good performance (Collins, 2002). We can use Passive-Aggressive (Crammer et al., 2006) or other ideas from large-margin training, computing the step size by dividing a non-negative **loss** $\ell(\mathbf{y}_i, \hat{\mathbf{y}})$ by the squared norm of the difference in the feature vectors, $\|\mathbf{f}(\mathbf{y}_i, \mathbf{w}_i) - \mathbf{f}(\hat{\mathbf{y}}, \mathbf{w}_i)\|^2$. A reasonable choice of loss function is the Hamming loss, which is the number of incorrect tag predictions (Taskar et al., 2003; Tsochantaridis et al., 2004).

9.4 Conditional random fields

Structured perceptron works well in practice, but sometimes we need probabilities $p(\mathbf{y} \mid \mathbf{w})$, and structured perceptron doesn't give us that. To fill this gap, the Conditional Random Field (CRF) is a probabilistic conditional model for sequence labeling; just as structured perceptron is built on the perceptron classifier, conditional random fields are built on the logistic regression classifier.

$$p(\mathbf{y} \mid \mathbf{x}) = \frac{e^{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{x})}}{\sum_{\mathbf{y}' \in \mathcal{Y}} e^{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}', \mathbf{x})}} \quad (9.46)$$

We can again move to structured prediction,

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$p(\mathbf{y} \mid \mathbf{w}) = \frac{e^{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w})}}{\sum_{\mathbf{y}' \in \mathcal{T}(\mathbf{w})} e^{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}', \mathbf{w})}}. \quad (9.47)$$

This is called a Conditional Random Field, because it models the sequence labeling task as a Markov random field, and estimates the probability of a set of variables **conditioned** on the others (as opposed to jointly, which the HMM does). We will need the same restriction on the scoring function as in Structured Perceptron: $\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_m \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$.

An important note about CRFs is that the joint probability $p(\mathbf{w}_{1:M}, \mathbf{y}_{1:M})$ is simply the unnormalized conditional probability:

$$p(\mathbf{y} \mid \mathbf{w}) = \frac{p(\mathbf{y}, \mathbf{w})}{\sum_{\mathbf{y}'} p(\mathbf{y}', \mathbf{w})} \quad (9.48)$$

$$p(\mathbf{y}, \mathbf{w}) = e^{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w})} \quad (9.49)$$

Aside: Maximum Entropy Markov Models (MEMMS)*

Suppose we define

$$p(\mathbf{y} \mid \mathbf{w}) = \prod_{m=1}^M p(y_m \mid \mathbf{w}_{1:M}, y_{1:m-1}) \quad (9.50)$$

$$\approx \prod_{m=1}^M p(y_m \mid \mathbf{w}_{1:M}, y_{m-1}). \quad (9.51)$$

We can then define each local probability $p(y_m \mid \mathbf{w}_{1:M}, y_{m-1})$ as a logistic regression model,

$$p(y_m \mid \mathbf{w}_{1:M}, y_{m-1}) = \frac{\exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, y_m, y_{m-1}))}{\sum_{y' \in \mathcal{T}} \exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_{1:M}, y', y_{m-1}))}. \quad (9.52)$$

Recall that logistic regression is sometimes called **maximum entropy**, and observe that we are making a Markov assumption. Thus the name **Maximum Entropy Markov Model**.

Inference in the MEMM can be performed with the Viterbi algorithm, choosing

$$v_m(k) = \max_{k'} P(Y_m = k \mid Y_{m-1} = k', \mathbf{w}_{1:M}) v_{m-1}(k'). \quad (9.53)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

The local decision model $p(y_m \mid \mathbf{w}_{1:M}, y_{m-1})$ can be trained as a standard logistic regression classifier. The problem with this model is that learning to optimize individual tagging decisions is not the same as learning to produce optimal tag sequences. The local classifier is trained with the true value of y_{m-1} , not the value likely to be produced by the classifier — so, not necessarily the value that we are most likely to see in a test set tagging situation. This introduces a problem that Lafferty et al. (2001) call **label bias**. Put another way, the MEMM allows structured **prediction**, but it does not perform structured **learning**.

Conditional random fields are designed to solve this problem: they elegantly unify structured prediction with a structured learning framework.

Decoding in CRFs

Decoding — predicting $\hat{\mathbf{y}}$ that maximizes $p(\mathbf{y} \mid \mathbf{w})$ — can be performed with the Viterbi algorithm. The key observation is that the decoding problem does not depend on the denominator of $p(\mathbf{y} \mid \mathbf{w})$,

$$\begin{aligned}
 \hat{\mathbf{y}} &= \arg \max_{\mathbf{y}} \log p(\mathbf{y} \mid \mathbf{w}) \\
 &= \arg \max_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w}) - \log \sum_{\mathbf{y}' \in \mathcal{T}(\mathbf{w})} e^{\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}', \mathbf{w})} \\
 &= \arg \max_{\mathbf{y}} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{y}, \mathbf{w}) \\
 &= \arg \max_{\mathbf{y}} \boldsymbol{\theta}^\top \sum_{m=1}^{M+1} \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \\
 &= \arg \max_{\mathbf{y}} \sum_{m=1}^{M+1} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m).
 \end{aligned}$$

This is identical to the decoding equation from structured perceptron.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Learning in CRFs

Learning in CRFs is a little more complicated. As with logistic regression, we need to learn weights to minimize the regularized negative log conditional probability,

$$\begin{aligned}\ell &= \sum_{i=1}^N -\log p(\mathbf{y}_i \mid \mathbf{w}_i; \boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|^2, \\ &= -\sum_i \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_i, \mathbf{y}_i) + \log \sum_{\mathbf{y}' \in \mathcal{T}(\mathbf{w}_i)} \exp(\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}_i, \mathbf{y}')) + \lambda \|\boldsymbol{\theta}\|^2,\end{aligned}$$

where λ controls the amount of regularization. As in logistic regression, the gradient includes is a difference between observed and expected counts:

$$\begin{aligned}\frac{d\ell}{d\theta_j} &= \sum_i \text{count}(\mathbf{w}_i, \mathbf{y}_i)_j - E_{\mathbf{y} \mid \mathbf{w}_i; \boldsymbol{\theta}}[\text{count}(\mathbf{w}_i, \mathbf{y})_j] + \lambda \theta_j \\ \text{count}(\mathbf{w}_i, \mathbf{y}_i)_j &= \sum_m^M f_j(\mathbf{w}_i, y_{i,m}, y_{i,m-1}, m)\end{aligned}$$

For example:

- If feature j is $\langle T : CC, DT \rangle$, then $\text{count}(\mathbf{w}_i, \mathbf{y}_i)_j$ is the count of times DT follows CC in the sequence \mathbf{y}_i .
- If feature j is $\langle M : -thy, JJ \rangle$, then $\text{count}(\mathbf{w}_i, \mathbf{y}_i)_j$ is the count of words ending in *-thy* in \mathbf{w}_i that are tagged JJ in \mathbf{y}_i .

The expected feature counts are more difficult to compute.

$$E_{\mathbf{y} \mid \mathbf{w}; \boldsymbol{\theta}}[\text{count}(\mathbf{w}_i, \mathbf{y})_j] = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w}_i)} P(\mathbf{y} \mid \mathbf{w}_i; \boldsymbol{\theta}) f_j(\mathbf{w}_i, \mathbf{y}) \quad (9.54)$$

This looks bad: we have to sum over an exponential number of labelings again. But remember that the feature function decomposition implies,

$$f_j(\mathbf{w}, \mathbf{y}) = \sum_m f_j(\mathbf{w}, y_m, y_{m-1}, m). \quad (9.55)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

This means we can compute the expectation as,

$$E_{\mathbf{y}|\mathbf{w};\theta}[\text{count}(\mathbf{w}, \mathbf{y})_j] = \sum_{\mathbf{y} \in \mathcal{T}(\mathbf{w})} p(\mathbf{y} | \mathbf{w}; \theta) f_j(\mathbf{w}, \mathbf{y}) \quad (9.56)$$

$$= \sum_{\mathbf{y} \in \mathcal{T}(\mathbf{w})} p(\mathbf{y} | \mathbf{w}; \theta) \sum_m^M f_j(\mathbf{w}, y_m, y_{m-1}, m) \quad (9.57)$$

$$= \sum_m^M \sum_{\mathbf{y} \in \mathcal{T}(\mathbf{w})} p(\mathbf{y} | \mathbf{w}; \theta) f_j(\mathbf{w}, y_m, y_{m-1}, m) \quad (9.58)$$

$$= \sum_m^M \sum_{k, k' \in \mathcal{T}} \sum_{\mathbf{y} \in \mathcal{T}(\mathbf{w}): Y_{m-1}=k, Y_m=k'} p(\mathbf{y} | \mathbf{w}; \theta) f_j(\mathbf{w}, k', k, m) \quad (9.59)$$

$$= \sum_m^M \sum_{k, k' \in \mathcal{T}} f_j(\mathbf{w}, k', k, m) \sum_{\mathbf{y} \in \mathcal{T}(\mathbf{w}): y_{m-1}=k, y_m=k'} p(\mathbf{y} | \mathbf{w}; \theta) \quad (9.60)$$

$$= \sum_m^M \sum_{k, k' \in \mathcal{T}} f_j(\mathbf{w}, k, k', m) Pr(Y_{m-1} = k, Y_m = k' | \mathbf{w}_{1:M}; \theta) \quad (9.61)$$

Consider $Pr(Y_{m-1} = k, Y_m = k' | \mathbf{w}_{1:M}; \theta)$: this is the probability of traversing the trellis edge $\langle m-1, k \rangle \rightarrow \langle m, k' \rangle$, conditioned on the entire observation $\mathbf{w}_{1:M}$. We can compute the tag bigram marginal as,

$$Pr(Y_{m-1} = k, Y_m = k' | \mathbf{w}_{1:M}; \theta) = Pr(Y_{m-1} = k, Y_m = k' | \mathbf{w}_{1:m-1}, w_m, \mathbf{w}_{m+1:M}; \theta) \quad (9.62)$$

$$= Pr(Y_{m-1} = k, \mathbf{w}_{1:m-1}) Pr(Y_m = k', w_m | Y_{m-1}=k) P(\mathbf{w}_{m+1:M} | Y_m = k') \quad (9.63)$$

In the HMM, we already know how to compute the first term, using the forward algorithm. The third term corresponds to the **backward probability**, which is computed from an equivalent backward algorithm. In this section, we will define versions of these algorithms that use log-linear probability distributions rather than generative probabilities as in the HMM. Given these probabilities, the expected feature counts can be computed efficiently if we know the **marginal** probabilities $P(Y_m = k', Y_{m-1} = k | \mathbf{w}; \theta)$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

The Forward-backward algorithm in CRFs

To compute the gradient of the CRF objective with respect to the weights θ , we need marginal probabilities of tag bigrams,

$$P(Y_m = k', Y_{m-1} = k | \mathbf{w}_{1:M}) = \frac{P(Y_m = k', Y_{m-1} = k, \mathbf{w}_{1:M})}{p(\mathbf{w}_{1:M})}. \quad (9.64)$$

To compute these, recall that in the CRF,

$$p(\mathbf{y} | \mathbf{w}) = \frac{\exp(\theta^\top \mathbf{f}(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}'} \exp(\theta^\top \mathbf{f}(\mathbf{w}, \mathbf{y}'))} \quad (9.65)$$

$$\propto \exp(\theta^\top \mathbf{f}(\mathbf{w}, \mathbf{y})) \quad (9.66)$$

$$= \exp\left(\theta^\top \sum_m \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)\right) \quad (9.67)$$

$$= \prod_m \exp(\theta^\top \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)) \quad (9.68)$$

$$= \prod_m \psi_m(y_m, y_{m-1}), \quad (9.69)$$

where $\psi_m(k, k') \triangleq \exp(\theta^\top \mathbf{f}(\mathbf{w}, Y_m = k, Y_{m-1} = k', m))$. This quantity is sometimes called a **potential**. A Hidden Markov Model (HMM) can also be expressed in terms of a product of potential functions, $\psi_m(k, k') = p(w_m | y_m) \times p(y_m | y_{m-1})$.

Using this expression for the probability of a label sequence, we will compute the desired marginal probabilities in terms of three quantities:

$$Z \triangleq \sum_{\mathbf{y}'} \prod_m \psi_m(y'_m, y'_{m-1}) \quad (9.70)$$

$$\alpha_m(k) \triangleq \sum_{\mathbf{y}'_{1:m-1}} \psi_m(k, y_{m-1}) \prod_n^{m-1} \psi_n(y'_n, y'_{n-1}) \quad (9.71)$$

$$\beta_m(k) \triangleq \sum_{\mathbf{y}'_{m+1:M}} \psi_{m+1}(y_{m+1}, k) \prod_{n=m+2}^M \psi_n(y'_n, y'_{n-1}), \quad (9.72)$$

where Z is the **normalization constant** (also called the partition function), α is the set of forward variables (just like in the forward algorithm defined earlier), and β is the set of backward variables.

(c) Jacob Eisenstein 2014-2015. Work in progress.

If we assume that the first tag and last tag have special values, $y_0 = \langle \text{START} \rangle$ and $y_M = \langle \text{STOP} \rangle$, then by definition,

$$Z = \alpha_M(\langle \text{STOP} \rangle) \quad (9.73)$$

$$= \beta_0(\langle \text{START} \rangle) \quad (9.74)$$

$$= \sum_{\mathbf{y}} p(\mathbf{w}, \mathbf{y}) \quad (9.75)$$

$$= p(\mathbf{w}) \quad (9.76)$$

The forward variables

We define the first sum as $\alpha_m(k)$, called the **forward variable**. This is just like in the forward algorithm defined earlier for Hidden Markov Models; again, we can derive a recursive expression for α as follows:

$$\alpha_0(\langle \text{START} \rangle) \triangleq 1 \quad (9.77)$$

$$\alpha_m(k) \triangleq \sum_{\mathbf{y}'_{1:m-1}} \psi_m(k, y'_{m-1}) \prod_{n=1}^{m-1} \psi_n(y'_n, y'_{n-1}) \quad (9.78)$$

$$= \sum_{Y_{m-1}=k'} \psi_m(k, k') \sum_{\mathbf{y}'_{1:m-2}} \psi_{m-1}(k', y'_{m-2}) \prod_{n=1}^{m-2} \psi_n(y'_n, y'_{n-1}) \quad (9.79)$$

$$= \sum_{Y_{m-1}=k'} \psi_m(k, k') \alpha_{m-1}(k'), \quad (9.80)$$

where we simply substitute in the recursive expression for $\alpha_{m-1}(k')$. This variable is computed from left to right, with each α_m depending on α_{m-1} .

The backward variables

We define the second sum as $\beta_m(k)$, called the **backward variable**. We can derive a recursive expression for β as follows:

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$\beta_M(\langle \text{STOP} \rangle) \triangleq 1 \quad (9.81)$$

$$\beta_m(k) \triangleq \sum_{\mathbf{y}'_{m+1:M}} \psi_{m+1}(y_{m+1}, k) \prod_{n=m+2}^M \psi_n(y'_n, y'_{n-1}), \quad (9.82)$$

$$= \sum_{Y_{m+1}=k'} \psi_{m+1}(k', k) \sum_{\mathbf{y}_{m+2:M}} \psi_{m+2}(y_{m+2}, k') \prod_{n=m+3}^M \psi_n(y_n, y_{n-1}) \quad (9.83)$$

$$= \sum_{Y_{m+1}=k'} \psi_{m+1}(k', k) \beta_{m+1}(k') \quad (9.84)$$

This variable is computed from right to left, with each β_m depending on β_{m+1} .

Computing the marginals Now, the product $\alpha_m(k)\beta_m(k)$ is equal to

$$\alpha_m(k)\beta_m(k) = \sum_{\mathbf{y}_{1:m-1}} \left(\prod_{n=1}^{m-1} \psi_n(y_n, y_{n-1}) \right) \psi_m(k, y_{m-1}) \quad (9.85)$$

$$\times \sum_{\mathbf{y}_{m+1:M}} \psi_{m+1}(y_{m+1}, k) \left(\prod_{n=m+2}^M \psi_n(y_n, y_{n-1}) \right) \quad (9.86)$$

$$= \sum_{\mathbf{y}_{1:M}: Y_m=k} \prod_{n=1}^M \psi_n(y_n, y_{n-1}), \quad (9.87)$$

which is exactly equal to the sum of the unnormalized probabilities of all sequences in which $Y_m = k$. To compute the normalized probability, we simply divide by Z , so that,

$$P(Y_m = k \mid \mathbf{w}_{1:M}) = \frac{\alpha_m(k)\beta_m(k)}{Z} \quad (9.88)$$

To obtain the unnormalized probability of a tag-to-tag transition $\langle Y_m = k, Y_{m-1} =$

(c) Jacob Eisenstein 2014-2015. Work in progress.

k'), we compute,

$$\alpha_{m-1}(k')\psi_m(k, k')\beta_m(k) = \sum_{\mathbf{y}_{1:m-2}} \left(\prod_{n=1}^{m-2} \psi_n(y_n, y_{n-1}) \right) \psi_{m-1}(k', y_{m-2}) \quad (9.89)$$

$$\times \psi_m(k, k') \sum_{\mathbf{y}_{m+1:M}} \psi_{m+1}(y_{m+1}, k) \prod_{n=m+2}^M \psi_n(y_n, y_{n-1}) \quad (9.90)$$

$$= \sum_{\mathbf{y}_{1:M}: Y_m=k, Y_{m-1}=k'} \prod_{n=1}^M \psi_n(y_n, y_{n-1}), \quad (9.91)$$

which is exactly equal to the unnormalized probability of all sequences in which $Y_m = k$ and $Y_{m-1} = k'$. To compute the normalized probability, we again divide by Z , so that,

$$P(Y_m = k, Y_{m-1} = k' \mid \mathbf{w}_{1:M}) = \frac{\alpha_{m-1}(k')\psi_m(k, k')\beta_m(k)}{Z} \quad (9.92)$$

Learning in CRFs: summary

The overall procedure looks like logistic regression, in which we compute a gradient from the difference between observed and expected counts.

- We use forward-backward to compute expected feature counts under the conditional probability distribution $p(\mathbf{y}_{1:M} \mid \mathbf{w}_{1:M}; \boldsymbol{\theta})$:

$$E[f_j(\mathbf{y}, \mathbf{w})] = \sum_m \sum_{k, k'} Pr(Y_m = k, Y_{m-1} = k' \mid \mathbf{w}_{1:M}) f_j(\mathbf{w}, k, k', m) \quad (9.93)$$

$$= \frac{1}{Z} \sum_m \sum_{k, k'} \alpha_m(k)\psi_m(k, k')\beta_{m+1}(k') f_j(\mathbf{w}, k, k', m). \quad (9.94)$$

- The gradient of the log-likelihood is the difference between feature counts and expected counts, $\mathbf{f}(\mathbf{y}, \mathbf{w}) - E[\mathbf{f}(\mathbf{y}, \mathbf{w})]$.
- We then update $\boldsymbol{\theta}$, typically by using either quasi-newton batch optimization or online stochastic gradient descent. The `CRFsuite` package implements several learning algorithms (<http://www.chokkan.org/software/crfsuite/>).

We iterate this procedure until convergence.

9.5 Unsupervised sequence labeling*

In unsupervised sequence labeling, we want to induce a Hidden Markov Model from a corpus of unannotated text $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_N$; this is an example of the general problem of **structure induction**, which is the unsupervised version of **structure prediction**. The tags that result from unsupervised sequence labeling might be useful for some downstream task, or for better understanding the language's inherent structure; or, we might want to do probability density estimation for sequences, as in gesture or activity recognition (Mittra and Acharya, 2007). Another reason would be to do semi-supervised learning, imputing tag sequences for unlabeled data. For part-of-speech tagging, often we use a tag dictionary which lists the allowed tags for each word, simplifying the problem (Christodoulopoulos et al., 2010).

In any case, we can perform unsupervised learning by applying expectation-maximization (EM). In the M-step, we compute the HMM parameters from expected counts:

$$\begin{aligned} Pr_E(W = i \mid Y = k) &= \phi_{k,i} = \frac{E[\text{count}(W = i, Y = k)]}{E[\text{count}(Y = k)]} \\ Pr_T(Y_m = k \mid Y_{m-1} = k') &= \lambda_{k',k} = \frac{E[\text{count}(Y_m = k, Y_{m-1} = k')]}{E[\text{count}(Y_{m-1} = k')]} \end{aligned}$$

The expected counts are computed in the E-step, using the forward and backward variables.

$$E[\text{count}(W = i, Y = k)] = \sum_m P(Y_m = k \mid \mathbf{w}) \delta(W_m = i) \quad (9.95)$$

$$= \sum_m \frac{P(Y_m = k, \mathbf{w}_{1:m}) p(\mathbf{w}_{m+1:M} \mid Y_m = k)}{p(\mathbf{w}_{1:M})} \delta(w_m = i) \quad (9.96)$$

$$= \frac{1}{\alpha_M(\langle \text{STOP} \rangle)} \sum_m \alpha_m(k) \beta_m(k) \delta(w_m = i) \quad (9.97)$$

We use the chain rule to separate $\mathbf{w}_{1:m}$ and $\mathbf{w}_{m+1:M}$, and then use the definitions of the forward and backward variables. In the final step, we normalize by $p(\mathbf{w}_{1:M}) = \alpha_M(\langle \text{STOP} \rangle) = \beta_0(\langle \text{START} \rangle)$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$E[\text{count}(Y_m = k, Y_{m-1} = k')] = \sum_m P(Y_m = k, Y_{m-1} = k' \mid \mathbf{w}) \quad (9.98)$$

$$\begin{aligned} &\propto \sum_m P(Y_{m-1} = k', \mathbf{w}_{1:m-1}) P(w_{m+1:M} \mid Y_m = k) \\ &\quad \times P(w_m, Y_m = k \mid Y_{m-1} = k') \end{aligned} \quad (9.99)$$

$$\begin{aligned} &= \sum_m P(Y_{m-1} = k', \mathbf{w}_{1:m-1}) P(w_{m+1:M} \mid Y_m = k) \\ &\quad \times p(w_m \mid Y_m = k) P(Y_m = k \mid Y_{m-1} = k') \end{aligned} \quad (9.100)$$

$$= \sum_m \alpha_{m-1}(k') \beta_m(k) \phi_{k,w_m} \lambda_{k' \rightarrow k} \quad (9.101)$$

Again, we use the chain rule to separate out $\mathbf{w}_{1:m-1}$ and $\mathbf{w}_{m+1:M}$, and use the definitions of the forward and backward variables. The final computation also includes the parameters ϕ and λ , which govern (respectively) the emission and transition properties between w_m, y_m , and y_{m-1} . Note that the derivation only shows how to compute this to a constant of proportionality; we would divide by $p(\mathbf{w}_{1:M})$ to go from the joint probability $P(Y_{m-1} = k', Y_m = k, \mathbf{w}_{1:M})$ to the desired conditional $P(Y_{m-1} = k', Y_m = k \mid \mathbf{w}_{1:M})$. As in the CRF, the joint probability $p(\mathbf{w}_{1:M})$ is given by the forward variable $\alpha_M(\langle \text{STOP} \rangle)$.

Linear dynamical systems

The forward-backward algorithm can be viewed as Bayesian state estimation in a discrete state space. In a continuous state-space, $y_m \in \mathbb{R}$, the equivalent algorithm is the **Kalman Smoother**. It also computes marginals $p(y_m \mid \mathbf{x}_{1:M})$, using a similar two-step algorithm of forward and backward passes. Instead of computing a table of values at each step ($\alpha_m(k)$ and $\beta_m(k)$), we would compute a probability density function $q_{y_m}(y_m; \mu_m, \Sigma_m)$, characterized by a mean μ_m and a covariance Σ_m around the latent state. Connections between the Kalman Smoother and the forward-backward algorithm are elucidated by Minka (1999) and Murphy (2012).

Alternative unsupervised learning methods

As noted in section 4.4, expectation-maximization is just one of many techniques for structure induction. One alternative is to use a family of randomized algorithms called **Markov Chain Monte Carlo (MCMC)**. In these algorithms, we compute a marginal distribution over the latent variable \mathbf{y} **empirically**, by drawing

random samples. The randomness explains the “Monte Carlo” part of the name; typically, we employ a Markov Chain sampling procedure, meaning that each sample is drawn from a distribution that depends only on the previous sample (and not on the entire sampling history). A simple MCMC algorithm is **Gibbs Sampling**, in which we iteratively sample each y_m conditioned on all the others (Finkel et al., 2005):

$$p(y_m \mid \mathbf{y}_{-m}, \mathbf{w}_{1:M}) \propto p(w_m \mid y_m) p(y_m \mid \mathbf{y}_{-m}). \quad (9.102)$$

Gibbs Sampling has been applied to unsupervised part-of-speech tagging by Goldwater and Griffiths (2007). *Beam sampling* is a more sophisticated sampling algorithm, which randomly draws entire sequences $\mathbf{y}_{1:M}$, rather than individual tags y_m ; this algorithm was applied to unsupervised part-of-speech tagging by Van Gael et al. (2009).

EM is guaranteed to find only a local optimum; MCMC algorithms will converge to the true posterior distribution $p(\mathbf{y}_{1:M} \mid \mathbf{w}_{1:M})$, but this is only guaranteed in the limit of infinite samples. Recent work has explored the use of **spectral learning** for latent variable models, which use matrix and tensor decompositions to provide guaranteed convergence under mild assumptions (Song et al., 2010; Hsu et al., 2012).

Chapter 10

Context-free grammars

So far we've explored finite-state models, which correspond to regular languages.

- **representations:** (weighted) finite state automata
- **probabilistic models:** HMMs (as a special case), CRFs
- **algorithms:** Viterbi, Forward-Backward, $\mathcal{O}(NK^2)$ time complexity.
- **linguistic phenomena:**
 - morphology
 - language models
 - part-of-speech disambiguation
 - named entity recognition (chunking)

Is the finite state representation enough for natural language?

10.1 Is English a regular language?

Regular languages are closed under intersection:

- $K \cap L$ is the set of strings in both K and L
- $K \cap L$ is regular iff K and L are regular

How to prove English is not regular:

- Let K be the set of grammatical English sentences
- Let L be some regular language

- Show that the intersection is not regular

We're going to prove this using center embedding:

1. *The cat is fat.*
2. *The cat that the dog chased is fat.*
3. **The cat that the dog is fat.*
4. *The cat that the dog that the monkey kissed chased is fat.*
5. **The cat that the dog that the monkey chased is fat.*

Proof sketch:

- K is the set of grammatical english sentences.
It excludes sentences (3) and (5).
- L is the regular language *the cat (that N)⁺ V_t ⁺ is fat.*
- The language $L \cap K$ is *the cat (that N)ⁿ V_t ⁿ is fat.*

It is important to understand that the issue here is not just infinite repetition or productivity; FSAs can handle productive phenomena like *the big red smelly plastic figurine*. It is specifically the center-embedding phenomenon, because this leads to the same structure as the classic $a^n b^n$ language. What do you think of this argument?

Is deep center embedding really part of English?

Karlsson (2007) searched for multiple (phrasal) center embeddings in corpora from 7 languages:

- Very few examples of double embedding
- Only 13 examples of triple embedding (none in speech)
- Zero examples of quadruple embeddings

Note that we can build an FSA to accept center-embedding up to any finite depth. Chomsky and many linguists distinguish between

- **Competence:** the fundamental abilities of the (idealized) human language processing system

(c) Jacob Eisenstein 2014-2015. Work in progress.

- **Performance:** real utterances produced by speakers, subject to non-linguistic factors such as cognitive limitations

Even if English *as performed* is regular, the underlying generative grammar may be context-free... **or beyond**. There is a similar proof that at least some languages are not context-free! I'll post slides with this proof idea.

How much expressiveness do we need?

- Shieber (1985) makes a similar argument, showing that case agreement in Swiss-German cross-serial constructions is homomorphic to a formal language $wa^mb^nc^md^ny$, which is weakly non-context free. In response to the objection that all attested constructions are finite, Shieber writes:

Down this path lies tyranny. Acceptance of this argument opens the way to proofs of natural languages as regular, nay, **finite**.

- In practice, many real constructions are much simpler to handle in context-free rather than finite-state representations:

*The **processor** has 10 million times fewer transistors on it than today's typical microprocessors, **runs** much more slowly, and **operates** at five times the voltage...*

- The easy way:

$$\begin{aligned} S &\rightarrow \text{NN VP} \\ \text{VP} &\rightarrow \text{VP3S} \mid \text{VPN3S} \mid \dots \\ \text{VP3S} &\rightarrow \text{VP3S, VP3S, and VP3S} \mid \text{VBZ} \mid \text{VBZ NP} \mid \dots \end{aligned}$$

- The hard way: build an FST that basically replicates all of English grammar for VPs with 3S and non-3S subjects.
- Mainstream parsing focuses on CFGs, but there is some work on “mildly” context-sensitive grammars.

10.2 Context-Free Languages

The Chomsky Hierarchy Every automaton defines a language, and different classes of automata define different classes of languages. The Chomsky Hierarchy:

- **finite-state automata** define **regular** languages;
- **pushdown automata** define **context-free** languages;
- **Turing machines** define **recursively-enumerable** languages.

In the Chomsky hierarchy, context-free languages (CFLs) are a strict generalization of regular languages.

regular	context-free
regular expressions	context-free grammars (CFGs)
finite-state machines	pushdown automata
paths	derivations

Context-free grammars define CFLs. They are sets of permissible *productions* which allow you to **derive** strings composed of surface symbols.

$$\begin{aligned}
 S &\rightarrow NP VP_1 \\
 NP &\rightarrow the\ N \mid NP\ RELCLAUSE \\
 RELCLAUSE &\rightarrow that\ NP\ V_t \\
 V_t &\rightarrow ate \mid chased \mid befriended \mid \dots \\
 N &\rightarrow cat \mid dog \mid monkey \mid \dots \\
 VP_1 &\rightarrow is\ fat
 \end{aligned}$$

An important feature of CFGs is *recursion*, in which a nonterminal can be derived from itself.

More formally: a CFG is a tuple $\langle N, \Sigma, R, S \rangle$:

- N a set of non-terminals
- Σ a set of terminals (distinct from N)
- R a set of productions, each of the form
 $A \rightarrow \beta$, where $A \in N$ and $\beta \in (\Sigma \cup N)^*$
- S a designated start symbol

- Context free grammars provide rules for generating strings.
 - RHS: a non-terminal $\in N$
 - LHS: a sequence of terminals or non-terminals, $\{n, \sigma\}^*, n \in N, \sigma \in \Sigma$.

- A **derivation** t is a sequence of steps from S to a surface string $w \in \Sigma^*$, which is the yield of the derivation.
- If $\exists t : s \in \text{yield}(t)$, then w is grammatical. Equivalently, for grammar G , we say $|\mathcal{T}_G(w)| \geq 1$.
- If $\exists t, t' : w \in \text{yield}(t) \wedge s \in \text{yield}(t')$, then w is ambiguous. Equivalently, for grammar G , we say $|\mathcal{T}_G(w)| > 1$.
- A derivation can be viewed as trees or as bracketings, as shown in Figure 10.1.

Semantics Ideally, each derivation will have a distinct semantic interpretation, and all possible interpretations will be represented in some derivation.

$$(\text{NP}(\text{NP Ban } (\text{PP on } (\text{NP nude dancing })))$$

$$(\text{PP on } (\text{NP Governor's desk })))$$

$$(\text{NP Ban } (\text{PP on } (\text{NP}(\text{NP nude dancing })$$

$$(\text{PP on } (\text{NP Governor's desk })))))$$

Sadly, this is not always the case.

$$(\text{NP}(\text{JJ nice }) (\text{JJ little }) (\text{NN car }))$$

$$(\text{NP}(\text{JJ nice }) (\text{NP}(\text{JJ little }) (\text{NN car })))$$

$$(\text{NP}(\text{JJ nice }) (\text{NP}(\text{JJ little }) (\text{NP}(\text{NN car }))))$$

10.3 Constituents

- In natural language grammars, the non-terminals should reflect syntactic categories.
- Bracketed substrings (e.g., *sushi with chopsticks*) are called **constituents**.
- There are several tests for constituency, including:
 - substitution
 - coordination
 - movement

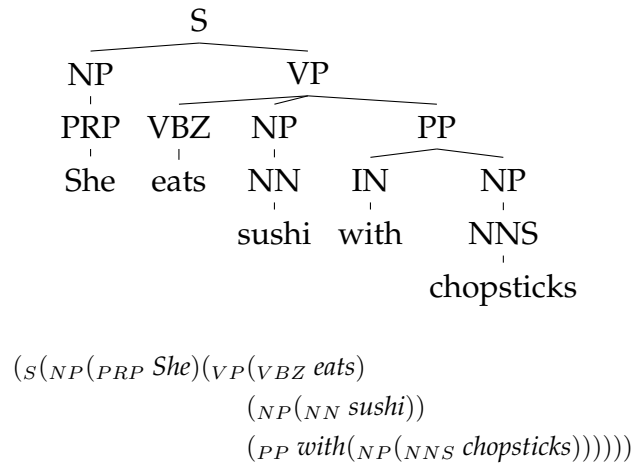
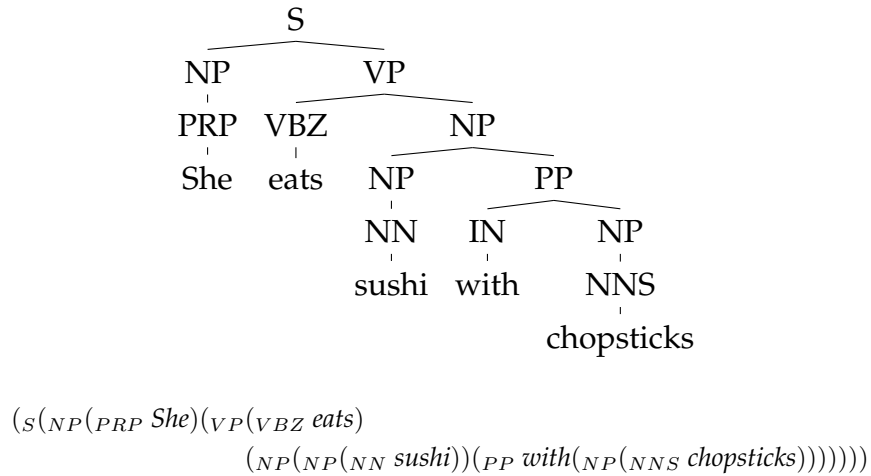


Figure 10.1: Two derivations of the same sentence, shown as both parse trees and bracketings

Substitution Constituents generated by the same non-terminal should be substitutable in many contexts:

- $(_{NP}\text{ The ban })$ is on the desk.
- $(_{NP}\text{ The Governor's desk })$ is on the desk.
- $(_{NP}\text{ The ban on dancing on the desk })$ is on the desk.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- $*(_{PP} \text{ On the desk }) \text{ is on the desk.}$

A more precise test for whether a set of substrings constitute a single category is whether they can be replaced by the same pronouns.

- $(_{NP} \text{ It }) \text{ is on the desk.}$

What about verbs?

- $I (_{V} \text{ gave }) \text{ it to Anne.}$
- $I (_{V} \text{ taught }) \text{ it to Anne.}$
- $I (_{V} \text{ gave }) \text{ Anne a fish}$
- $*I (_{V} \text{ taught }) \text{ Anne a fish}$

This suggests we need nonterminals which distinguish verbs based on the arguments they can take. The technical name for this is *subcategorization*.

Coordination Constituents generated by the same non-terminal can usually be *coordinated* using words like *and* and *or*:

- $\text{We fought } (_{PP} \text{ on the hills }) \text{ and } (_{PP} \text{ in the hedges }).$
- $\text{We fought } (_{ADVP} \text{ as well as we could }).$
- $*\text{We fought } (_{ADVP} \text{ as well as we could }) \text{ and } (_{PP} \text{ in the hedges }).$

This too doesn't always work:

- $\text{She } (_{VP} \text{ went }) (_{PP} \text{ to the store }).$
- $\text{She } (_{VP} \text{ came }) (_{PP} \text{ from the store }).$
- $\text{She } (\text{ went to }) \text{ and } (\text{ came from }) \text{ the store.}$

Movement Valid constituents can be moved as a unit, preserving grammaticality.

- Passivization
 - $(\text{The governor}) \text{ banned } (\text{nude dancing on his desk})$
 - $(\text{Nude dancing on his desk}) \text{ was banned by } (\text{the governor})$
- Wh- movement

- *(Nude dancing was banned) on (the desk).*
- *(The desk) is where (nude dancing was banned)*
- Topicalization
 - *(He banned nude dancing) to appeal to conservatives.*
 - *To appeal to conservatives, (he banned nude dancing).*

10.4 A simple grammar of English

Noun phrases

Let's start with noun phrases:

- *She sleeps* (Pronoun)
- *Arlo sleeps* (Proper noun)
- *Fish sleep* (Mass noun)
- *The fish sleeps* (determiner + noun)
- *The blue fish sleeps* (DT + JJ + NN)
- *The girl from Omaha sleeps* (NP + PP)
- *The student who ate 15 donuts sleeps* (NP + RelClause)

So overall, we can summarize this fragment as

$$\begin{aligned} \text{NP} &\rightarrow \text{PRP} \mid \text{NNP} \mid \text{DT NOM} \\ \text{NOM} &\rightarrow \text{ADJP NOM} \mid \text{NN} \\ \text{NP} &\rightarrow \text{NP PP} \mid \text{NP RELCLAUSE} \end{aligned}$$

We're leaving out some detail, like pluralization and possessives, but you get the idea.

Adjectival and prepositional phrases

- *Very funny*
- *The large, blue fish*
- *The man from la Mancha*

$$\begin{aligned}\text{ADJP} &\rightarrow \text{JJ} \mid \text{RB ADJP} \mid \text{JJ ADJP} \\ \text{PP} &\rightarrow \text{IN NP} \mid \text{TO NP}\end{aligned}$$

Verb phrases

- *She sleeps*
- *She sleeps restlessly*
- *She sleeps at home*
- *She eats sushi*¹
- *She gives John sushi*

$$\text{VP} \rightarrow \text{V} \mid \text{VP RB} \mid \text{VP PP} \mid \text{V NP} \mid \text{V NP NP} \mid \text{V NP RB}$$

But what about **She sleeps sushi* or **She speaks John Japanese*?

- Classes of verbs can take different numbers of arguments.
- This is called **subcategorization**

$$\begin{aligned}\text{VP} &\rightarrow \text{V-INTRANS} \mid \text{V-TRANS NP} \mid \text{V-DITRANS NP NP} \\ \text{VP} &\rightarrow \text{VP RB} \mid \text{VP PP}\end{aligned}$$

We would also need to handle modal and auxiliary verbs that allow us to create complex tenses, like *She will have eaten sushi* but not **She will have eats sushi*.

Sentences

- *She eats sushi*

$$\text{S} \rightarrow \text{NP VP}$$

- *Sometimes, she eats sushi*

$$\text{S} \rightarrow \text{ADVP S}$$

- *In Japan, she eats sushi*

$$\text{S} \rightarrow \text{PP S}$$

¹Sushi examples from Julia Hockenmaier

- What about **I eats sushi*, **She eat sushi??*

$$S \rightarrow \text{NP.3S VP.3S} \mid \text{NP.N3S VP.N3S}$$

In general, we need **features** to capture this kind of agreement.

Conjunctions

- *She eats sushi and candy*

$$\text{NP} \rightarrow \text{NP and NP}$$

- *She eats sushi and drinks soda*

$$\text{VP} \rightarrow \text{VP and VP}$$

- *She eats sushi and he drinks soda*

$$S \rightarrow S \text{ and } S$$

- *fresh and tasty sushi*

$$\text{ADJP} \rightarrow \text{JJ and JJ}$$

We'd need a little more cleverness to properly cover groups larger than two.

Odds and ends

- *I gave sushi to the girl **who eats sushi***. This is a relative clause,

$$\text{RELCLAUSE} \rightarrow \text{who VP} \mid \text{that VP}$$

- *I took sushi from the man **offering sushi***. This is a gerundive postmodifier.

$$\text{NOM} \rightarrow \text{NOM GERUNDVP}$$

$$\text{GERUNDVP} \rightarrow \text{VBZ} \mid \text{VBZ NP} \mid \text{VBZ PP} \mid \dots$$

- ***Can** she eat sushi?* (notice it's not *eats*)

$$S \rightarrow \text{AUX NP VP}$$

- ... and many more

(c) Jacob Eisenstein 2014-2015. Work in progress.

10.5 Grammar design

Our goal is a grammar that avoids

- **Overgeneration:** deriving strings that are not grammatical.
- **Undergeneration:** failing to derive strings that are grammatical.

To avoid undergeneration, we would need thousands of productions.

Typically, grammars are defined in conjunction with large-scale **treebank** annotation projects.

- An annotation guideline specifies the non-terminals and how they go together.
- The annotators then apply these guidelines to data.
- The grammar rules can then be read off the data.

The Penn Treebank (PTB) contains one million parsed words of Wall Street Journal text (Marcus et al., 1993).

10.6 Grammar equivalence and normal form

There may be many grammars that express the same context-free language:

- Grammars are **weakly equivalent** if they generate the same strings.
- Grammars are **strongly equivalent** if they generate the same strings **and** assign the same phrase structure to each string.

In Chomsky Normal Form (CNF), all productions are either:

$$A \rightarrow BC$$

$$A \rightarrow a$$

All CFGs can be converted into a CNF grammar that is weakly equivalent — meaning that it generates exactly the same set of strings. As we will soon see, this conversion is very useful for parsing algorithms.

In CNF, all productions have either two or zero non-terminals on the right-hand side. To deal with productions that have more than two non-terminals on the RHS, we create new “dummy” non-terminals. For example, if we have $W \rightarrow X Y Z$, we can replace this production with two productions: $X \rightarrow W X \setminus W$ and

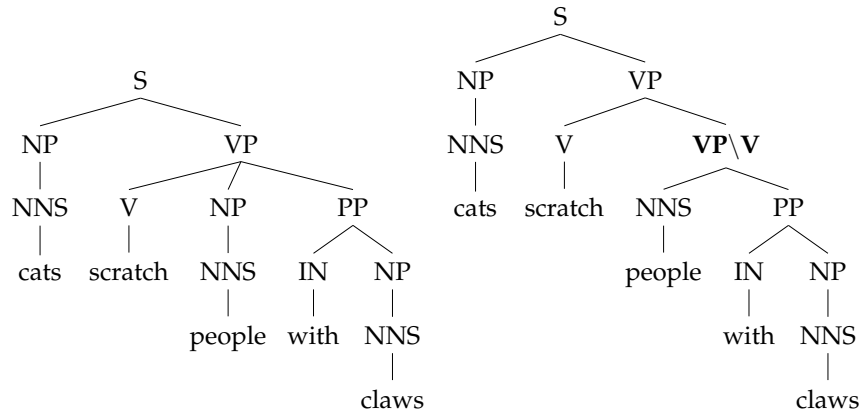


Figure 10.2: Binarization of the $VP \rightarrow V NP PP$ production

$X \setminus W \rightarrow Y Z$, where $X \setminus W$ is a new dummy non-terminal. Figure 10.2 conveys this idea in a real example.

Note that *people with claws* was not a constituent in the original grammar, but it is a constituent in the binarized grammar. Therefore, after parsing it is important to take care to “un-binarize” the resulting parse.

What about unary productions, such as $NP \rightarrow NNS$? While we could easily deal with this in the grammar, as we will see, in practice it is best dealt with by modifying the parsing algorithm itself.

[**todo: recap**]

Chapter 11

CFG Parsing

Parsing is the process of determining whether a sentence is in a context-free language, by searching for a legal derivation. Here are some possible approaches:

- **Top-down:** start with the start symbol, and see if we can derive the sentence.
- **Bottom-up:** combine the observed symbols using whatever productions we can, until we reach the start symbol
- **Left-to-right:** move through the input, incrementally building a parse tree

Before we get into these different possibilities, let's see whether exhaustive search is possible. Suppose we only have one non-terminal, X , and it has binary productions

$$X \rightarrow X X$$

$$X \rightarrow \textit{the girl} \mid \textit{ate sushi} \mid \dots$$

How many different ways could we derive a sentence in this language? This is equal to the number of binary bracketings of the words in the sentence, which is a Catalan number. Catalan numbers grow **super-exponentially** in the length of the sentence, $C_n = \frac{(2n)!}{(n+1)!n!}$. As with sequence labeling, we cannot search the space of possible derivations naïvely; we will again rely on dynamic programming to search efficiently by reusing shared substructures.

11.1 CKY parsing

CKY is a bottom-up parsing allows us to test whether a sentence is in a context-free language, without considering all possible parses. First we form small con-

stituents, then try to merge them into larger constituents.

Let's start with an example grammar:

$$\begin{aligned} S &\rightarrow VP \ NP \\ NP &\rightarrow NP \ PP \mid we \mid sushi \mid chopsticks \\ PP &\rightarrow P \ NP \\ P &\rightarrow with \\ VP &\rightarrow VP \ NP \mid VP \ PP \mid eat \end{aligned}$$

Suppose we encounter the sentence *We eat sushi with chopsticks*.

- The first thing that we notice is that we can apply unary productions to obtain $NP \ VP \ NP \ P \ NP$
- Next, we can apply a binary production to merge the first $NP \ VP$ into an S .
- Or we could merge $VP \ NP$ into VP
- ... and so on

Let's systematize this. The CKY algorithm (Algorithm 2) incrementally constructs a table t , where each cell $t[i, j]$ contains the set of nonterminals that can derive the span $w_{i:j}$. If $S \in t[0, N]$, then w is in the language defined by the grammar.

Algorithm 2 The CKY algorithm for CFG parsing

```

1: for  $j \leftarrow 1, N$  do
2:    $t[j, j - 1] \leftarrow \{X \mid X \rightarrow w_j \in R\}$ 
3:   for  $i \leftarrow j - 2, 0$  do
4:     for  $k \leftarrow i + 1, j - 1$  do
5:        $t[i, j] \leftarrow t[i, j] \cup \{X \mid X \rightarrow YZ \in R, Y \in t[i, k], Z \in t[k, j]\}$ 

```

The CKY algorithm assumes that all productions with non-terminals on the RHS are binary. What do we do when this is not true?

- For productions with more than two elements on the right-hand side, we binarize, creating additional non-terminals. For example, if we have the production $VP \rightarrow VNPNP$ (for ditransitive verbs), we might convert to $VP \rightarrow VP_{ditrans}/NPNP$, and then add the production $VP_{ditrans}/NP \rightarrow VNP$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- What about unary productions like $S \rightarrow VP \rightarrow V \rightarrow eat$? To handle this case, we compute the *unary closure* of each non-terminal. For example, if the grammar includes $S \rightarrow VP$ and $VP \rightarrow V$, then we add $S \rightarrow V$ to the unary closure of S . Then for each entry $t[i, j]$ in the table, for each non-terminal $A \in t[i, j]$, we add all elements of the reflexive unary closure for A to $t[i, j]$.

Complexity What is the complexity of CKY?

- Space complexity: $\mathcal{O}(M^2|N|)$
- Time complexity: $\mathcal{O}(M^3|R|)$
- M is length of sentence,
 $|N|$ is the number of non-terminals,
 $|R|$ is the number of production rules
- But in practice... It's worse than worst-case! (Figure 11.1)
- Because longer sentences “unlock” more of the grammar.

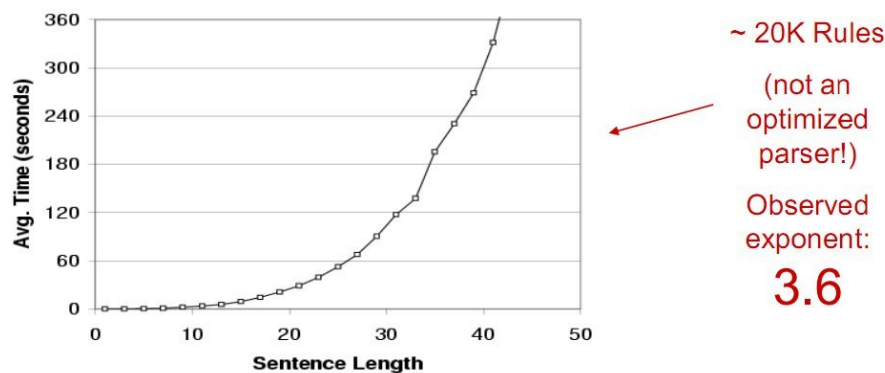


Figure 11.1: Figure from Dan Klein’s lecture slides

11.2 Ambiguity in parsing

- Syntactic ambiguity is endemic to natural language:¹

¹Examples borrowed from Dan Klein

- Attachment ambiguity: *we eat sushi with chopsticks, I shot an elephant in my pajamas.*
 - Modifier scope: *southern food store*
 - Particle versus preposition: *The puppy tore up the staircase.*
 - Complement structure: *The tourists objected to the guide that they couldn't hear.*
 - Coordination scope: *"I see," said the blind man, as he picked up the hammer and saw.*
 - Multiple gap constructions: *The chicken is ready to eat*
- In morphology, we didn't just want to know which derivational forms are *legal*, we wanted to know which were *likely*.
 - Syntactic parsing is all about choosing among the many, many legal parses for a given sentence.

Here's another example, which we've seen before:

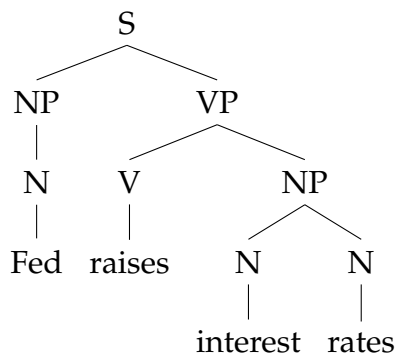


Figure 11.2: An example parse

- A minimal grammar permits 36 parses!
- Real-size broad coverage grammars permit millions of parses of typical sentences.
- Collins (2013) gives another example, *She announced a program to promote safety in trucks.*

Classical parsers faced a tradeoff:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- broad coverage with tons of ambiguity...
- or limited coverage in exchange for constraints on ambiguity

Consequently, deterministic parsers produced no analysis for many sentences.

Local solutions

Some ambiguity can be resolved locally:

- [*imposed* [*a ban* [*on asbestos*]]]
- [*imposed* [*a ban*]] [*on asbestos*]]
- Hindle and Rooth (1990) proposed a likelihood ratio test:

$$LR(v, n, p) = \frac{p(p | v)}{p(p | n)} = \frac{p(on | imposed)}{p(on | ban)}$$

where we select VERB attachment if $LR(v, n, p) > 1$.

- But the likelihood-ratio approach ignores important information, like the phrase being attached.
 - ...[*it* [*would end* [*its venture* [*with Maserati*]]]]
 - ...[*it* [*would end* [*its venture*]] [*with Maserati*]]
- The likelihood ratio gets this wrong
 - $p(with | end) = \frac{607}{5156} = 0.118$
 - $p(with | venture) = \frac{155}{1442} = 0.107$

Other features (e.g., *Maserati*) argue for noun attachment. How can we add them?

Machine learning solutions Ratnaparkhi et al. (1994) propose a classification-based approach, using logistic regression (maximum-entropy):

$$P(\text{Noun attachment} | \text{would end its venture with Maserati}) = \frac{e^{\theta^\top \mathbf{f}(\text{would end its venture with Maserati})}}{1 + e^{\theta^\top \mathbf{f}(\text{would end its venture with Maserati})}}$$

Features include n-grams and word classes from hierarchical word clustering; accuracy is roughly 80%.

Collins and Brooks (1995) argued that attachment depends on four **heads**:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- the preposition (*with*)
- the VP attachment site (*end*)
- the NP attachment site (*venture*)
- the NP to be attached (*Maserati*)

They propose a backoff-based approach:

- First, look for counts of the tuple $\langle with, Maserati, end, venture \rangle$
- If none, try $\langle with, Maserati, end \rangle + \langle with, end, venture \rangle + \langle with, Maserati, venture \rangle$
- If none, try $\langle with, Maserati \rangle + \langle with, end \rangle + \langle with, venture \rangle$
- If none, try $\langle with \rangle$

Accuracy is roughly 84%. This approach of combining relative frequency estimation, smoothing, and backoff was very characteristic of late 1990s statistical NLP.

Beyond local solutions

Framing the problem as attachment ambiguity is limiting:

- assumes the parse is mostly done, leaving just a few attachment ambiguities to solve
- But realistic sentences have more than a few syntactic interpretations.
- Attachment decisions are interdependent:
 - *Cats scratch people with claws with knives.*
 - We may want to attach *with claws* to *scratch*.
 - But then we have nowhere to put *with knives*.

The task of statistical parsing is to produce a single analysis that resolves all syntactic ambiguities.

S	→ NP VP	0.9
S	→ S CC S	0.1
NP	→ N	0.2
NP	→ DT N	0.3
NP	→ N NP	0.2
NP	→ JJ NP	0.2
NP	→ NP PP	0.1
VP	→ V	0.4
VP	→ V NP	0.3
VP	→ V NP NP	0.1
VP	→ VP PP	0.2
PP	→ P NP	1.0

Table 11.1: A fragment of an example probabilistic context-free grammar (PCFG)

11.3 Statistical parsing with PCFGs

We want the parse τ that maximizes $p(\tau \mid \mathbf{w})$.

$$\begin{aligned}
 \arg \max_{\tau} p(\tau \mid \mathbf{w}) &= \arg \max_{\tau} \frac{p(\tau, \mathbf{w})}{p(\mathbf{w})} \\
 &= \arg \max_{\tau} p(\tau, \mathbf{w}) \\
 &= \arg \max_{\tau} p(\mathbf{w} \mid \tau) p(\tau) \\
 &= \arg \max_{\tau: \mathbf{w} = \text{yield}(\tau)} p(\tau)
 \end{aligned}$$

- The **yield** of a tree is the string of terminal symbols that can be read off the leaf nodes.
- The set $\{\tau : \mathbf{w} = \text{yield}(\tau)\}$ is exactly the set of all derivations of \mathbf{w} in a CFG G .

PCFGs extend the CFG by adding probability to each production, as shown in Table 11.1.

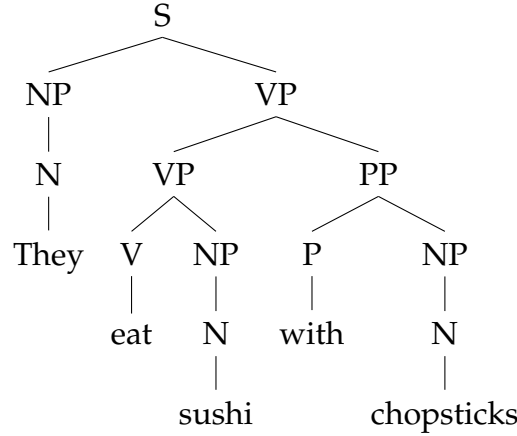


Figure 11.3: An alternative parse of our running example

The probabilities for all productions involving a single LHS must sum to 1,²

$$\sum_{\alpha} P(X \rightarrow \alpha) = 1.$$

The probability of a derivation $p(\tau)$ is just the product of all the productions. Consider the example in Figure 11.3, which we will score using the probabilities in Table 11.1.

The probability of this parse is:

$$p(\tau, w) = P(S \rightarrow NP VP) \tag{11.1}$$

$$\times P(NP \rightarrow N) \times P(N \rightarrow they) \tag{11.2}$$

$$\times P(VP \rightarrow VP PP) \tag{11.3}$$

$$\times P(VP \rightarrow V NP) \times P(V \rightarrow eat) \tag{11.4}$$

$$\times P(NP \rightarrow N) \times P(N \rightarrow sushi) \tag{11.5}$$

$$\times P(PP \rightarrow P NP) \times P(P \rightarrow with) \tag{11.6}$$

$$\times P(NP \rightarrow N) \times P(N \rightarrow chopsticks) \tag{11.7}$$

$$= 0.9 \times 0.2 \times 0.2 \times 0.3 \times 0.2 \times 1.0 \times 0.2 \tag{11.8}$$

$$\times \text{probability of terminal productions} \tag{11.9}$$

²Note that we use the P symbol here, because we are referring the probability of the **event** of the non-terminal X generating the right-hand side α .

Estimation

As in supervised HMMs, estimation is easy (for now!). We can estimate the production probabilities directly from a treebank, using relative frequency estimation. For example,

$$P(\text{VP} \rightarrow \text{VP PP}) = \frac{\text{count}(\text{VP} \rightarrow \text{VP PP})}{\text{count}(\text{VP})}$$

Three basic problems for PCFGs

Let $\tau \in T$ be a derivation, w be a sentence, and λ a PCFG.

- **Decoding:** Find $\hat{\tau} = \arg \max_{\tau} p(\tau, w; \lambda)$
- **Likelihood:** Find $p(w; \lambda) = \sum_{\tau} p(\tau, w; \lambda)$
- **(Unsupervised) Estimation:** Find $\arg \max_{\lambda} p(w_{1..N} | \lambda)$

These three problems are analogous to the problems identified by Rabiner (1989) for Hidden Markov Models.

	Sequences	Trees
model	HMM	PCFG
decoding	Viterbi algorithm	CKY
decoding complexity	$\mathcal{O}(M^2 K)$	$\mathcal{O}(M^3 R)$
likelihood	forward algorithm	inside algorithm
marginals	forward-backward	inside-outside

CKY with probabilities

It is not difficult to extend CKY to include probabilities or other weights. This is shown in Algorithm 3.

In the boolean semiring, we have $\oplus = \vee$ and $\otimes = \wedge$, so we recover the original CKY algorithm. In the probability semiring, we have $\oplus = \max$ and $\otimes = \times$.

$$t[Y, i, k] = P(Y \rightarrow w_{i:k}) \quad (11.10)$$

$$t[Z, k, j] = P(Z \rightarrow w_{k+1:j}) \quad (11.11)$$

$$t[X, i, j] = \max_{Y, Z, k} P(X \rightarrow Y Z) P(Y \rightarrow w_{i:k}) P(Z \rightarrow w_{k+1:j}) \quad (11.12)$$

$$(11.13)$$

Algorithm 3 CKY with weighted productions

```

for  $j \leftarrow 1, N$  do
  for all  $X \in \text{tags}(w_j)$  do
     $t[X, j, j - 1] \leftarrow P(X \rightarrow w_j)$ 
  for  $i \leftarrow (j - 2), 0$  do
    for all  $(X \rightarrow Y Z) \in R$  do
      for  $k \leftarrow (i + 1), (j - 1)$  do
         $t[X, i, j] \leftarrow t[X, i, j] \oplus (P(X \rightarrow Y Z) \otimes t[Y, i, k] \otimes t[Z, k + 1, j])$ 

```

The **inside algorithm** computes the probability of producing a span of text $w_{i:j}$ from a non-terminal X . To do this, we move to a semiring where $\oplus = +$,

$$t[X, i, j] = \sum_{Y, Z, k} P(X \rightarrow Y Z) P(Y \rightarrow w_{i:k}) P(Z \rightarrow w_{k+1:j}) \quad (11.14)$$

$$= P(X \rightarrow w_{i:j}). \quad (11.15)$$

11.4 Algorithms for PCFG Parsing

PCFGs score the probability of a derivation $P(\tau, w)$ as the product of all productions in τ

- CKY is a bottom-up algorithm for finding $\hat{\tau} = \arg \max_{\tau} p(\tau, w)$.
- The **inside algorithm** finds $p(w) = \sum_{\tau} p(\tau, w)$.
- These algorithms are related through semiring notation.
- These algorithms are bottom-up: they parse progressively larger spans until the entire sentence is parsed. This is efficient, but it is pretty implausible as a model of human parsing, since it seems unrelated to the way we hear and read language: left-to-right.

Shift-reduce is a left-to-right parsing algorithm, which you may find more cognitively plausible. It is related to the pushdown automata representation of context-free grammars: we move through the sentence while keeping a stack with infinite depth. At each step, we have two choices:

- **shift** the next word on to the stack;

(c) Jacob Eisenstein 2014-2015. Work in progress.

- **reduce** the stack by applying some production.

Each reduce move is a production in the derivation. If we can clear all the input and end up with just S on the stack, we have parsed the sentence correctly.

1. Initial state

Stack	Remaining Text
	the dog saw a man in the park

2. After one shift

Stack	Remaining Text
the	dog saw a man in the park

3. After reduce shift reduce

Stack	Remaining Text
<div>Det N</div> <div>the dog</div>	saw a man in the park

4. After recognizing the second NP

Stack	Remaining Text
	in the park

5. After building a complex NP

Stack	Remaining Text

6. Built a complete parse tree

Stack	Remaining Text
<pre> graph TD S --> NP1[NP] S --> VP[VP] NP1 --> Det1[Det] NP1 --> N1[N] Det1 --> the1[the] N1 --> dog[dog] VP --> V[V] VP --> NP2[NP] V --> saw[saw] NP2 --> Det2[Det] NP2 --> NP3[NP] Det2 --> a[a] NP3 --> N2[N] N2 --> man[man] NP3 --> P[P] NP3 --> NP4[NP] P --> in[in] NP4 --> Det3[Det] NP4 --> N3[N] Det3 --> the2[the] N3 --> park[park] </pre>	

Figure 11.4: Example of shift-reduce CFG parsing, from Bird et al. (2009)

How do we decide whether to shift or reduce?

- We could treat this as a classic search problem, and just backtrack when we get into trouble.
- Or, we could train a classifier to decide between shift and reduce (Ratnaparkhi, 1999; Yamada and Matsumoto, 2003). Note that we have a separate reduce action for each non-terminal in the grammar.

11.5 Parser evaluation

Before continuing to more advanced parsing algorithms, we need to consider how to measure parsing performance. Suppose we have a set of **reference parses** — the

ground truth — and a set of **system parses** that we would like to score. A simple solution would be **per-sentence accuracy**: the parser is scored by the proportion of sentences on which the system and reference parses exactly match.³ But we would like to assign *partial credit* for correctly matching parts of the reference parse. The PARSEval metrics do that, scoring each system parse via:

Precision, the fraction of brackets in the system parse that match a bracket in the reference parse.

Recall, the fraction of brackets in the reference parse that match a bracket in the system parse.

In labeled precision and recall, it is required to also match the non-terminals for each bracket; in unlabeled precision and recall, it is only required to match the bracketing structure. The F-measure is the harmonic mean of precision and recall.

In Figure 10.1, suppose the top tree is the system parse and the bottom tree is the reference parse. We have the following spans:

- $S \rightarrow w_{1:5}$: true positive
- $VP \rightarrow w_{2:5}$: true positive
- $NP \rightarrow w_{3:5}$: false positive
- $PP \rightarrow w_{4:5}$: true positive

So for this parse, we have a (labeled and unlabeled) precision of $\frac{3}{4} = 0.75$, and a recall of $\frac{3}{3} = 1.0$, for an F-measure of 0.86. The best automatic CFG parsers get an F-score of approximately 0.92 on the Penn Treebank (PTB) today (McClosky et al., 2006).

11.6 Improving PCFG parsing

Regardless of the parsing algorithm, pure PCFG parsing on Penn Treebank non-terminals (e.g., NP, VP) doesn't work well: a PCFG build from treebank probabilities scores $F = 0.72$. Why?

³Most parsing papers do not report results on this metric, but Finkel et al. (2008) find that a near-state-of-the-art parser finds the exact correct parse on 35% of sentences of length ≤ 40 , and on 62% of parses of length ≤ 15 in the Penn Treebank.

Problems with PCFG parsing

Substitutability Recall that substitutability is a criterion for constituency. Are NPs really substitutable? No, because some pronouns cannot be both subjects and objects (Figure 11.5).

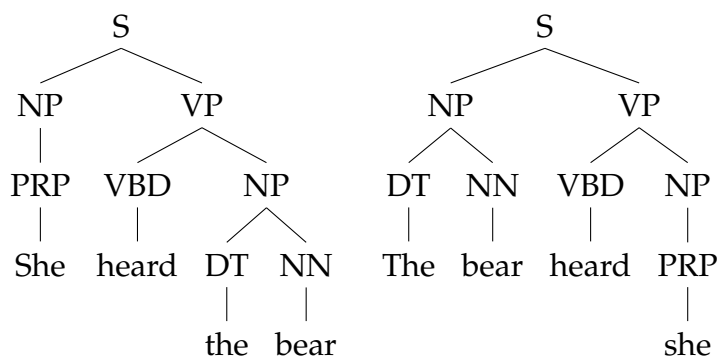


Figure 11.5: A grammar that allows *she* to take the object position is overgenerating.

We might address this problem by **splitting** the NP tag into nominative (*she*) and oblique (*her*) cases, but this distinction is only relevant for pronouns: other nouns can appear in either position.

A related point is that we have no flexibility on PP attachment. If $P(\text{NP} \rightarrow \text{NP PP}) > P(\text{VP} \rightarrow \text{VP PP})$, we will always prefer NP attachment; if not, we will always prefer VP attachment. More fine-grained NP and VP categories might allow us to make attachment decisions more accurately.

Semantic preferences In addition to grammatical constraints such as case marking, we have semantic preferences: for example, that conjoined entities should be similar (Figure 11.6).

Note that no PCFG can distinguish the parses in Figure 11.6! They contain exactly the same productions.

Subsumption There are several choices for annotating PP attachment

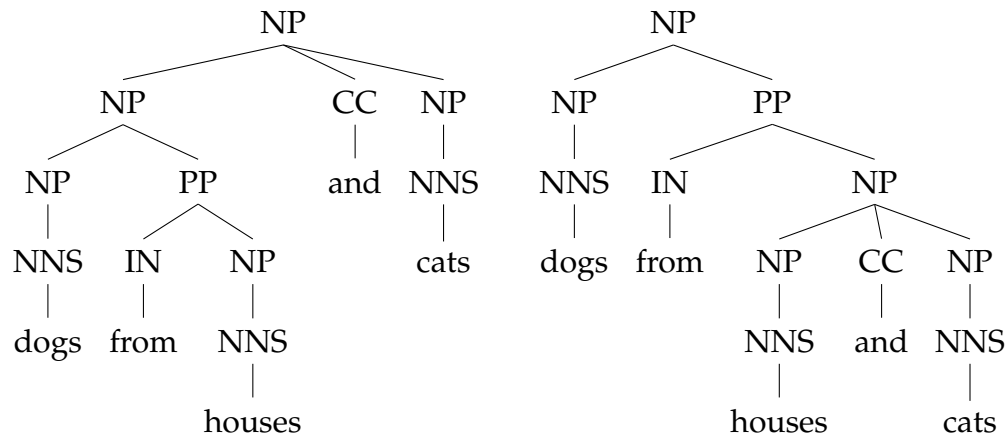
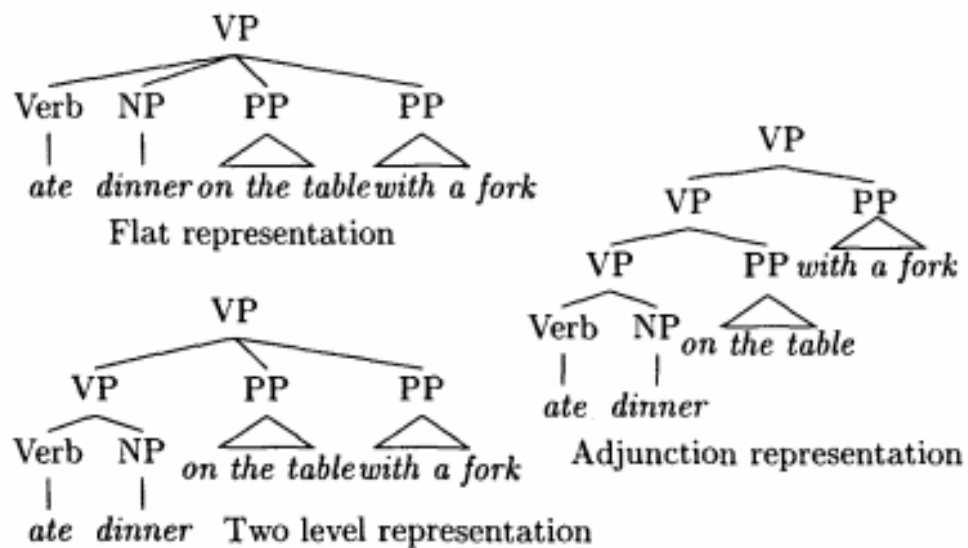


Figure 11.6: The first parse is arguably preferable because of the conjunction of phrases headed by *dogs* and *cats*. Example from Dan Klein's lecture slides.



Johnson (1998) shows that even though the two-level representation is chosen in the annotation, it can never be produced by a PCFG because the production is

subsumed.

$$\begin{aligned}P(\text{NP} \rightarrow \text{NP PP}) &= 0.112 \\P(\text{NP} \rightarrow \text{NP PP PP}) &= 0.006 \\P(\text{NP} \rightarrow \text{NP PP})P(\text{NP} \rightarrow \text{NP PP}) &= (0.112)^2 = 0.013\end{aligned}$$

The probability of applying the $\text{NP} \rightarrow \text{NP PP}$ production twice is greater than the probability of the two-PP production, so this production will never appear in a PCFG parse. Johnson shows that 9% of all productions are subsumed and can be removed from the grammar!

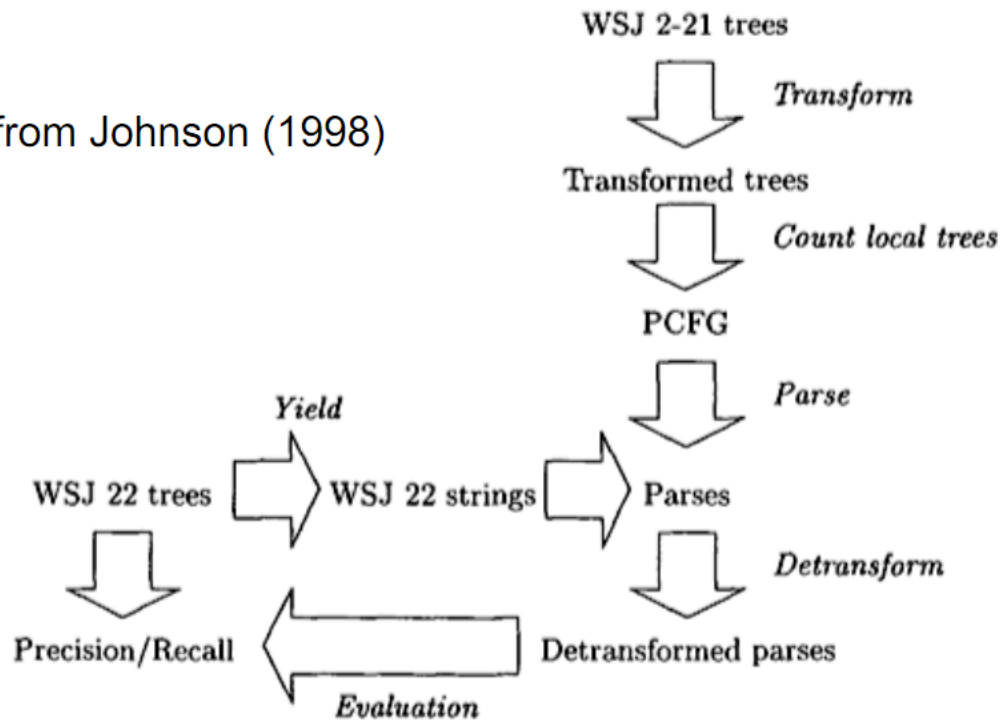
Refining PCFG parsing Modern generative parsing algorithms automatically refine these nonterminals in various ways.

- **Tree transformations:** automatically modify the parse trees — for example, Markovizing by labeling each non-terminal with its **parent**, as in NP-S.
- **Lexicalization:** label each non-terminal with its head **word**

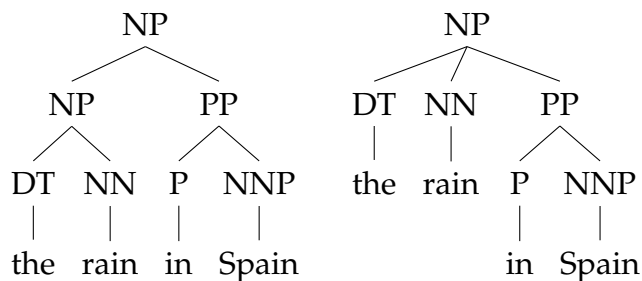
Tree transformations

Johnson proposed a series of transformations to PTB trees that improve parsing accuracy.

from Johnson (1998)



Flattening Johnson (1998) proposes “flattening” nested NPs to be more like VP structures.



Flattened rules are of course still context-free, but by reducing recursion, they allow more specific probabilities to be learned. This can eliminate the problems with rule subsumption that we saw earlier.

Parent annotation The expansion of an NP is highly dependent on its parent.

(c) Jacob Eisenstein 2014-2015. Work in progress.

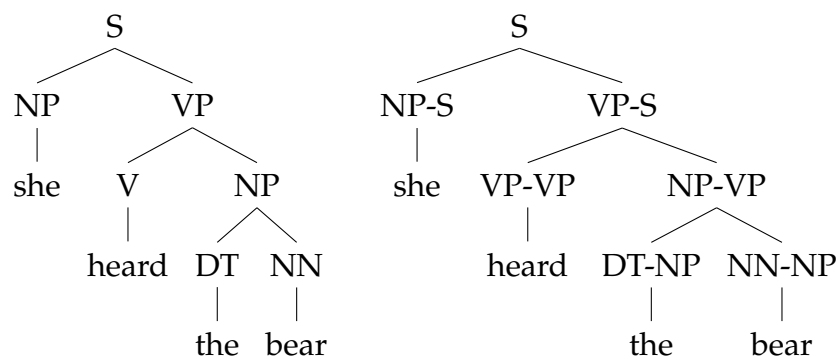
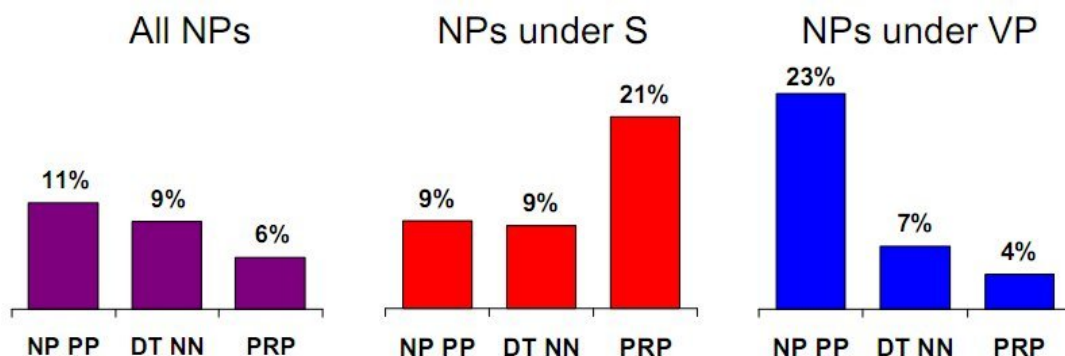


Figure 11.7: Parent annotation, sometimes called vertical Markovization (Klein and Manning, 2003).



$$\begin{aligned}
 P(\text{NP} \rightarrow \text{NP PP}) &= 11\% \\
 P(\text{NP-S} \rightarrow \text{NP PP}) &= 9\% \\
 P(\text{NP-VP} \rightarrow \text{NP PP}) &= 23\%
 \end{aligned}$$

PP adjunction is more likely for NPs that are under verb phrases (typically in object position) than for NPs that are under sentences (typically in subject position). We can capture this phenomenon via **parent annotation**: augmenting each non-terminal with the identity of its parent (Figure 11.7).

Parent annotation weakens the PCFG independence assumptions

- which could help accuracy by allowing more fine-grained distinctions
- or could hurt accuracy because of data sparseness

(c) Jacob Eisenstein 2014-2015. Work in progress.

Overall, the transformations proposed by Johnson (1998) improve performance on PTB parsing.

- Standard PCFG: 72% F-measure, 14,962 rules
- Parent-annotated PCFG: 80% F-measure, 22,773 rules
- In principle, parent annotation could have increased the grammar size much more dramatically, but many possible productions never occur, or are subsumed.

Lexicalization

A simple way to capture semantics is through the words themselves. We can annotate each non-terminal with **head** word of the phrase.

Head words are deterministically assigned according to a set of rules, sometimes called **head percolation rules**. In many cases, these rules are straightforward: the head of a NP \rightarrow DTN production is the noun, the head of a S \rightarrow NPVP production is the head of the VP, etc. But as always, there are a lot of special cases. Collins (2013) offers the following example for productions whose LHS is NP:

- **If** the RHS contains NN, NNS, or NNP, **then** choose the rightmost NN, NNS, or NNP.
- **Else If** the RHS contains an NP, **then** choose the leftmost NP
- **Else If** the rule contains a JJ, **then** choose the rightmost JJ (e.g., *Sandra is **the** best*)
- **Else If** the rule contains a CD (cardinal number), **then** choose the rightmost CD (e.g., *Marco is 27.*)
- **Else** choose the rightmost child.

A fragment of the head percolation rules used in many parsing systems are found in Table 11.2.⁴

The meaning of these rules is that to find the head of an S constituent, we first look for the rightmost VP child; if we don't find one, we look for the rightmost SBAR child, and so on down the list. Verb phrases are headed by left verbs (the head of *can walk home* is *walk*, since *can* is tagged MD), noun phrases are headed by the rightmost noun-like non-terminal (so the head of *the red cat* is *cat*), and

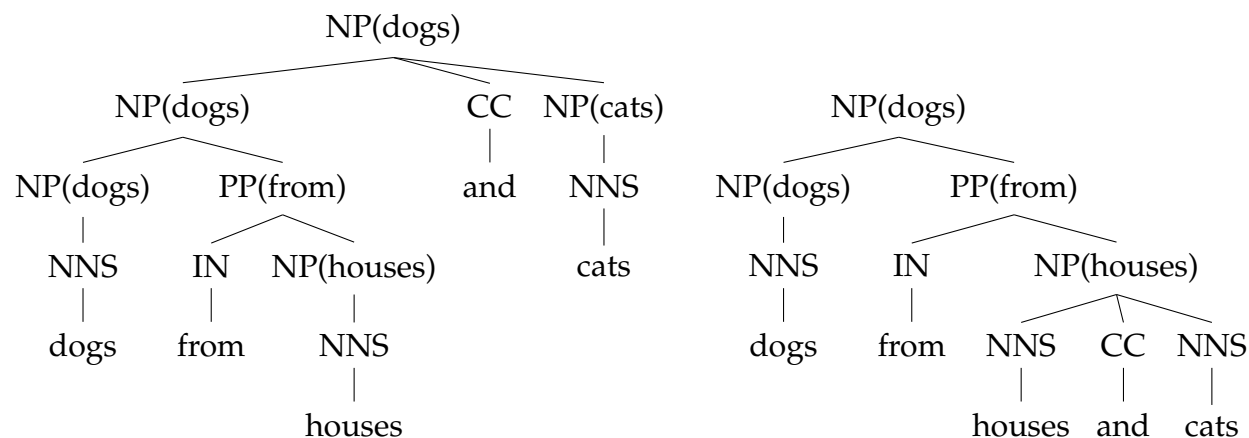
⁴From <http://www.cs.columbia.edu/~mcollins/papers/heads>

Non-terminal	Direction	Priority
S	right	VP SBAR ADJP UCP NP
VP	left	VBD VBN MD VBZ TO VB VP VBG VBP ADJP NP
NP	right	N* EX \$ CD QP PRP ...
PP	left	IN TO FW

Table 11.2: A fragment of head percolation rules

prepositional phrases are headed by the preposition (the head of *at Georgia Tech* is *at*). Some of these rules are somewhat arbitrary — there’s no particular reason why the head of *cats and dogs* should be *dogs* — but the point here is just to get some lexical information that can support parsing, not to make any deep claims about syntax.

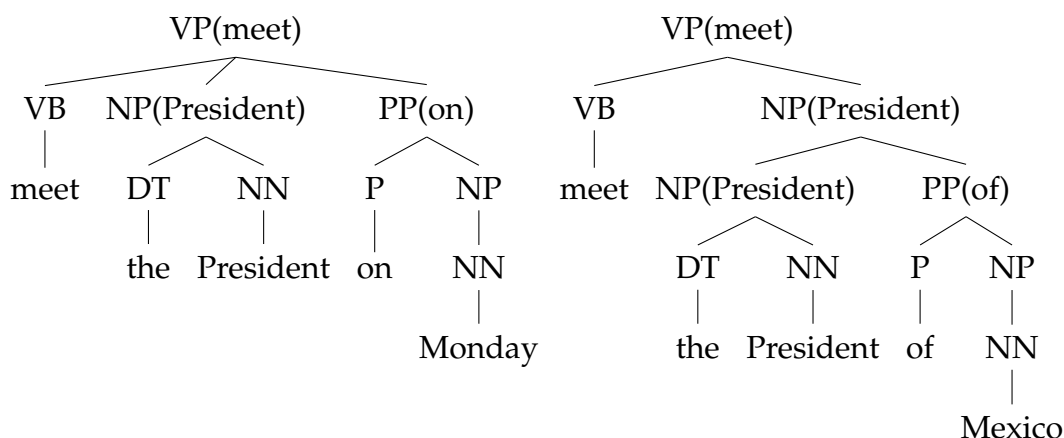
Given these rules, we can lexicalize the parse trees for some of our examples:



Example: coordination scope If $P(\text{NP} \rightarrow \text{NP}(\text{dogs}) \text{ CC } \text{NP}(\text{cats})) > P(\text{NP} \rightarrow \text{NP}(\text{houses}) \text{ CC } \text{NP}(\text{cats}))$, we should get the right parse.

Example: PP attachment

- $P(\text{VP}(\text{meet}) \rightarrow \alpha \text{ PP}(\text{on})) \gg P(\text{NP}(\text{President}) \rightarrow \beta \text{ PP}(\text{on}))$
- $P(\text{VP}(\text{meet}) \rightarrow \alpha \text{ PP}(\text{of})) \ll P(\text{NP}(\text{President}) \rightarrow \beta \text{ PP}(\text{of}))$
- In plain English: *Meeting* happens *on* things; *Presidents* are *of* things.



Subcategorization frames

$$\begin{aligned}
 P(\text{VP} \rightarrow \text{V NP NP}) &= 0.00151 \\
 P(\text{VP}(\text{said}) \rightarrow \text{V}(\text{said}) \text{ NP NP}) &= 0.00001 \\
 P(\text{VP}(\text{gave}) \rightarrow \text{V}(\text{gave}) \text{ NP NP}) &= 0.01980
 \end{aligned}$$

Lexicalization can capture fine-grained information that the Penn Treebank non-terminals ignore. This had a major impact on parsing accuracy, as shown in Table 11.3.

Vanilla PCFG	72%
Head-annotated PCFG (Johnson, 1998)	80%
Lexicalized PCFG (Collins, 1997, 2003; Charniak, 1997)	87-89%

Table 11.3: Penn Treebank parsing accuracies

Eugene Charniak: “To do better, it is necessary to condition probabilities on the actual words of the sentence. This makes the probabilities much tighter” (at a workshop at Johns Hopkins in 2000).

Algorithms for lexicalized parsing

Naively: we could just augment the non-terminals to include the cross-product of all PTB non-terminals and all words.

The number of possible productions: $\mathcal{O}(N^3V^2)$, where the size of the vocabulary V is $V \approx 10^5$, and N is the number of non-terminals. (The term on the

vocabulary size is V^2 , rather than V^3 , because the head of the entire parent must be identical to the head of one of the children.)

To perform lexicalized parsing, we can work bottom-up by building a table, similar to CKY. However, we need one additional piece of information: the location of the head word of each span. We should therefore store the elements $t[i, j, h, X]$, indicating a span from $w_{i:j}$, headed by w_h ($h \in i \dots j$), with parent node X .

To recursively construct $t[i, j, h, X]$, we need to consider two possibilities: either h is in the left child, or it is in the right child. If it is in the left child, then we have:

$$t_\ell[i, j, h, X] = \max_{s \geq h} \max_{m > s} \max_{X(w_h) \rightarrow Y(w_h)Z(w_m)} P(X(w_h) \rightarrow Y(w_h)Z(w_m)) \times t[i, s, h, Y] \times t[s, j, m, Z] \quad (11.16)$$

Otherwise, we have

$$t_r[i, j, h, X] = \max_{s < h} \max_{i < m \leq s} \max_{X(w_h) \rightarrow Y(w_m)Z(w_h)} P(X(w_h) \rightarrow Y(w_m)Z(w_h)) \times t[i, s, m, Y] \times t[s, j, h, Z] \quad (11.17)$$

We are building a table of size $\mathcal{O}(M^3 N)$. To fill in each cell we perform $\mathcal{O}(M^2 G)$ operations, taking maxes over two indices in the sentence, and over all rules. However, Eisner and Satta (1999) show that this time cost can be reduced back to $\mathcal{O}(M^3)$. A more serious problem is **estimation**: we must compute probabilities for $\mathcal{O}(N^3 V^2)$ productions. Charniak (1997) and Collins (1997, 2003) offer practical solutions, which decompose the production probabilities using various independence assumptions.

The Charniak Parser

The Charniak (1997) parser gives a relatively straightforward way to lexicalize PCFGs.

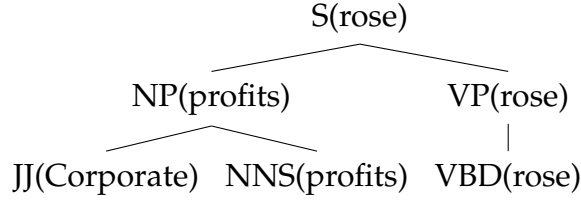
Recall that **head probabilities** capture “bilexical” phenomena, like the PP attachment in the example, *President of Mexico*. The idea of the Charniak (1997) parser is to represent the probability of each constituent by the product of two probabilities:

- The **rule probability**, $P(r \mid w_{\rho m}, t_m, t_{\rho m})$, where r is the rule, m is the index of the head of the left-hand side, t_m is the type of the left-hand side (non-terminal), and $t_{\rho m}$ is the type of the parent of m .

(c) Jacob Eisenstein 2014-2015. Work in progress.

- The **head probability**, $P(w_m \mid w_{\rho m}, t_m, t_{\rho m})$, where w_m is a head word.

Consider this example:



- The rule probability is

$$P(\text{NP} \rightarrow \text{JJ NNS} \mid w_m = \text{rose}, t_m = \text{NP}, t_{\rho m} = \text{S}) \quad (11.18)$$

- The head probability is

$$p(\text{profits} \mid w_{\rho m} = \text{rose}, t_m = \text{NP}, t_{\rho(m)} = \text{S}) \quad (11.19)$$

We would then multiply these probabilities to fill in the CKY table. This parser therefore combines two ideas that we have seen before:

- Head annotation, since both the rule and head probabilities depend on the parent type t and the grandparent type ℓ .
- Lexicalization, since the rule probability depends on the head word. Such rule probabilities can capture phenomena like verb complement frames:

Local Tree	come	take	think	want
VP \rightarrow V	9.5%	2.6%	4.6%	5.7%
VP \rightarrow V NP	1.1%	32.1%	0.2%	13.9%
VP \rightarrow V PP	34.5%	3.1%	7.1%	0.3%
VP \rightarrow V SBAR	6.6%	0.3%	73.0%	0.2%
VP \rightarrow V S	2.2%	1.3%	4.8%	70.8%
VP \rightarrow V NP S	0.1%	5.7%	0.0%	0.3%
VP \rightarrow V PRT NP	0.3%	5.8%	0.0%	0.0%
VP \rightarrow V PRT PP	6.1%	1.5%	0.2%	0.0%

- Bilexical probabilities, since the head probability depends on the head words of both the parent and child.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Estimating the Charniak parser The Charniak parser involves fewer parameters than a naive lexicalized PCFG. To estimate the relevant parameters in our example, we have

$$\begin{aligned}
 P_{\text{head}}(\text{profits} \mid t_m = \text{NP}, t_{\rho(m)} = \text{S}, w_{\rho(m)} = \text{rose}) \\
 &= \frac{\text{count}(w_m = \text{profits}, t_m = \text{NP}, t_{\rho(m)} = \text{S}, w_{\rho(m)} = \text{rose})}{\text{count}(t_m = \text{NP}, t_{\rho(m)} = \text{S}, w_{\rho(m)} = \text{rose})} \\
 P_{\text{rule}}(\text{NP} \rightarrow \text{JJ NNS} \mid w_{\rho(m)} = \text{rose}, t_m = \text{NP}, t_{\rho(m)} = \text{S}) \\
 &= \frac{\text{count}(\text{NP} \rightarrow \text{JJ NNS}, t_m = \text{NP}, t_{\rho(m)} = \text{S}, w_{\rho(m)} = \text{rose})}{\text{count}(t_m = \text{NP}, t_{\rho(m)} = \text{S}, w_{\rho(m)} = \text{rose})}
 \end{aligned}$$

The Penn Treebank provides is still the main dataset for syntactic analysis of English. Yet its 1M words is not nearly enough data to accurately estimate lexicalized models such as the Charniak parser, without smoothing.

For example, in 965K annotated constituent spans, there are

- 66 examples of WHADJP
- only 6 of these aren't *how much* or *how many*

In the example above (*corporate profits rose*), the unsmoothed head probability is zero, as estimated from the PTB: there are zero counts of *profits* headed by *rose* in the treebank (hard to believe, but that's what Charniak says).

In general, bilexical counts are going to be very sparse. But the “backed-off” probabilities give a reasonable approximation. We will interpolate between them.

Smoothing the Charniak Parser Head probability:

$$\begin{aligned}
 \hat{p}(w_m \mid t_m, w_{\rho(m)}, t_{\rho(m)}) &= \lambda_1 p_{\text{mle}}(w_m \mid t_m, w_{\rho(m)}, t_{\rho(m)}) \\
 &\quad + \lambda_2 p_{\text{mle}}(w_m \mid t_m, \text{cluster}(w_{\rho(m)}), t_{\rho(m)}) \\
 &\quad + \lambda_3 p_{\text{mle}}(w_m \mid t_m, t_{\rho(m)}) \\
 &\quad + \lambda_4 p_{\text{mle}}(w_m \mid t_m)
 \end{aligned}$$

, where $\text{cluster}(w_{\rho(m)})$ is the cluster of word $w_{\rho(m)}$, obtained by applying an automatic clustering method to **distributional** statistics (Pereira et al., 1993).

		$p(\text{profit} \mid \text{NP}, \text{rose}, \text{S})$	$P(\text{corp.} \mid \text{JJ}, \text{profit}, \text{NP})$
	$p(w_m \mid t_m, w_{\rho(m)}, t_{\rho(m)})$	0	.245
For example:	$p(w_m \mid t_m, c(w_{\rho(m)}), t_{\rho(m)})$.0035	.015
	$p(w_m \mid t_m, t_{\rho(m)})$.00063	.0053
	$p(w_m \mid t_m)$.00056	.0042

We have to tune $\lambda_1 \dots \lambda_4$, and an equivalent set of parameters for the rule probabilities.

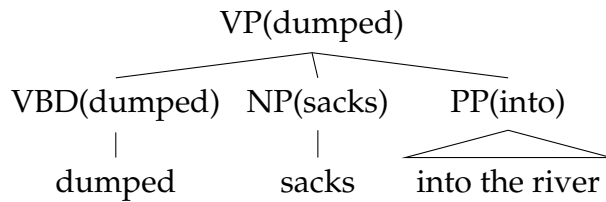
- The Charniak parser suffers from acute sparsity problems because it estimates the probability of entire rules.
- Another extreme would be to generate the children independently from each other.
e.g., $P(S \rightarrow NP VP) \approx P_L(S \rightarrow NP)P_R(S \rightarrow VP)$
- Collins (2003) and Charniak (2000) go for a compromise, conditioning on the parent and the head child.

The Collins Parser

- The Charniak parser focuses on lexical relationships between children and parents.
- The Collins (2003) parser focuses on relationships between adjacent children of the same parent. It decomposes each rule as,

$$X \rightarrow L_m L_{m-1} \dots L_1 H R_1 \dots R_{n-1} R_n$$

- Each L and R is a child constituent of X , and they are generated from the head H outwards.
- The outermost elements of L and R are special $\langle \text{STOP} \rangle$ symbols.



To model this rule, we would compute:

$$p(\text{VP}(\text{dumped}, \text{VBD}) \rightarrow [\langle \text{STOP} \rangle, \text{VBD}(\text{dumped}, \text{VBD}), \text{NP}(\text{sacks}, \text{NNS}), \text{PP}(\text{into}, \text{P}), \langle \text{STOP} \rangle])$$

- Here's the generative process:
 - Generate the head: $P(H \mid LHS) = P(\text{VBD}(\text{dumped}, \text{VBD}) \mid \text{VP}(\text{dumped}, \text{VBD}))$
 - Generate the left dependent: $P_L(\langle \text{STOP} \rangle \mid \text{VP}(\text{dumped}, \text{VBD}), \text{VBD}(\text{dumped}, \text{VBD}))$

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Generate the right dependent: $P_R(\text{NP}(\text{sacks}, \text{NNS}) \mid \text{VP}(\text{dumped}, \text{VBD}), \text{VBD}(\text{dumped}, \text{VBD}))$
 - Generate the right dependent: $P_R(\text{NP}(\text{into}, \text{PP}) \mid \text{VP}(\text{dumped}, \text{VBD}), \text{VBD}(\text{dumped}, \text{VBD}))$
 - Generate the right dependent: $P_R(\langle \text{STOP} \rangle \mid \text{VP}(\text{dumped}, \text{VBD}), \text{VBD}(\text{dumped}, \text{VBD}))$
- The rule probability is the product of these generative probabilities.
 - Collins parser also conditions on a “distance” of each constituent from the head.

Smoothing the Collins Parser Estimation is eased by factoring the rule probabilities, but smoothing is still needed.

$$\begin{aligned} \hat{P}(R_m(rw_m, rt_m) \mid \rho(m), hw, ht) = & \lambda_1 P_{mle}(R_m(rw_m, rt_m) \mid \rho(m), hw, ht) \\ & + \lambda_2 P_{mle}(R_m(rw_m, rt_m) \mid \rho(m), ht) \\ & + \lambda_3 P_{mle}(R_m(rw_m, rt_m) \mid \rho(m)) \end{aligned}$$

For example,

$$\begin{aligned} \hat{P}_R(\text{NP}(\text{sacks}, \text{NNS}) \mid \text{VP}(\text{dumped}, \text{VBD}), \text{dumped}, \text{VBD}) \\ = & \lambda_1 \hat{P}(\text{NP}(\text{sacks}, \text{NNS}) \mid \text{VP}, \text{dumped}, \text{VBD}) \\ & + \lambda_2 \hat{P}(\text{NP}(\text{sacks}, \text{NNS}) \mid \text{VP}, \text{VBD}) \\ & + \lambda_3 \hat{P}(\text{NP}(\text{sacks}, \text{NNS}) \mid \text{VP}) \end{aligned}$$

We set λ using Witten-Bell smoothing.

Bilexical dependencies The collins parser models bilexical dependencies in $P_{mle}(R_m(rw_m, rt_m) \mid \rho(m), hw, ht)$. Is it worth it?

Back-off level	Number of accesses	Percentage
0	3,257,309	1.49
1	24,294,084	11.0
2	191,527,387	87.4
Total	219,078,780	100.0

- In general, bilexical probabilities are rarely available...
- ...but they are active in 29% of the rules in **top-scoring** parses.
- Still, they don't seem to play a big role in accuracy (Bikel, 2004)

Summary of lexicalized parsing

Lexicalized parsing results in substantial accuracy gains:

Vanilla PCFG	72%
Parent-annotations (Johnson, 1998)	80%
(Charniak, 1997)	86%
(Collins, 2003)	87%

Table 11.4: Accuracies for lexicalized parsers

But the explosion in the size of the grammar required elaborate smoothing techniques and made parsing slow.

- Treebank syntactic categories are too coarse, but lexicalized categories may be too fine; modern approaches have sought middle ground.
- At the same time, natural language processing moved from generative models to more advanced machine learning techniques in the late 1990s and early 2000s, and researchers have worked to incorporate these techniques into parsing.

11.7 Modern constituent parsing

Reranking

Charniak and Johnson (2005) combine generative and discriminative models for parsing, using the idea of **reranking**. First, a generative model is used to identify its K -best parses. Then a discriminative ranker is trained to select the best of these parses. The discriminative model doesn't need to search over all parses — just the best K identified by the generative model. This means that it can use arbitrary features — such as structural features that capture parallelism and right-branching, which could not be easily incorporated into a bottom-up parsing model. This approach yields substantial improvements in accuracy on the Penn Treebank.

Refinement grammars

Klein and Manning (2003) revisit unlexicalized parsing, expanding on the ideas in (Johnson, 1998).

They apply two types of **Markovization**:

- Vertical Markovization, making the probability of each parsing rule depend not only on the type of the parent symbol, but also on its parent type. This is identical to the parent annotation proposed by Johnson (1998). The amount of vertical Markovization can be written v , with $v = 1$ indicating a standard PCFG.
- Horizontal Markovization, where the probability of each child depends on only some of its siblings. In a standard PCFG $h = \infty$, since there is no decomposition on the right-hand side of the rule. In the Collins parser, different settings of h were explored, with $h = 1$ indicating dependence only on the head, and $h = 2$ indicating dependence on the nearest sibling as well as the head.

A comparison of various Markovization parameters is shown in Figure 11.8:

Second, Klein and Manning note that the right level of linguistic detail is somewhere between treebank categories and individual words.

- Some parts-of-speech and non-terminals are truly substitutable: for example, *cat*/N and *dog*/N.
- But others are not: for example, *on*/PP behaves differently from *of*/PP. This is an example of **subcategorization**.

Vertical Order		Horizontal Markov Order				
		$h = 0$	$h = 1$	$h \leq 2$	$h = 2$	$h = \infty$
$v = 1$	No annotation	71.27 (854)	72.5 (3119)	73.46 (3863)	72.96 (6207)	72.62 (9657)
$v \leq 2$	Sel. Parents	74.75 (2285)	77.42 (6564)	77.77 (7619)	77.50 (11398)	76.91 (14247)
$v = 2$	All Parents	74.68 (2984)	77.42 (7312)	77.81 (8367)	77.50 (12132)	76.81 (14666)
$v \leq 3$	Sel. GParents	76.50 (4943)	78.59 (12374)	79.07 (13627)	78.97 (19545)	78.54 (20123)
$v = 3$	All GParents	76.74 (7797)	79.18 (15740)	79.74 (16994)	79.07 (22886)	78.72 (22002)

Figure 11.8: Performance for various Markovization levels (Klein and Manning, 2003).

- Similarly, Klein and Manning distinguish *and* and *but* from other coordinating conjunctions.
- They distinguish recursive and non-recursive noun phrases.

Figure 11.9 shows an example of an error that is corrected through the introduction of a new NP-TMP subcategory.

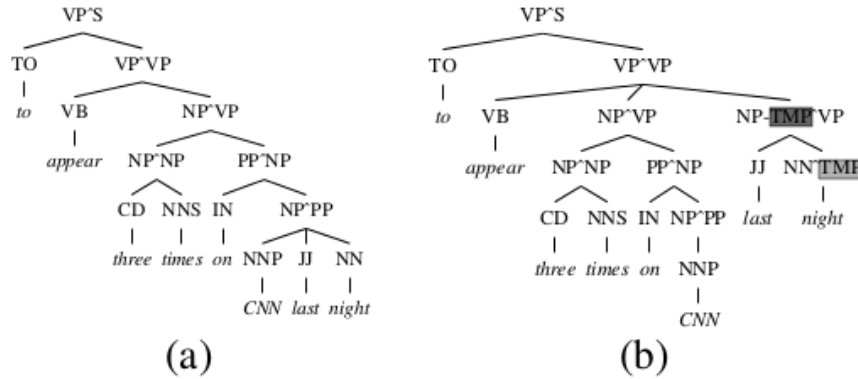


Figure 11.9: State-splitting creates a new non-terminal called NP-TMP, for temporal noun phrases. This corrects the PCFG parsing error in (a), resulting in the correct parse in (b).

(c) Jacob Eisenstein 2014-2015. Work in progress.

Klein and Manning (2003) analyze the improvements offered by each additional tag split, shown in Figure 11.10.

Annotation	Cumulative			Indiv.
	Size	F ₁	ΔF_1	ΔF_1
Baseline ($v \leq 2, h \leq 2$)	7619	77.77	—	—
UNARY-INTERNAL	8065	78.32	0.55	0.55
UNARY-DT	8066	78.48	0.71	0.17
UNARY-RB	8069	78.86	1.09	0.43
TAG-PA	8520	80.62	2.85	2.52
SPLIT-IN	8541	81.19	3.42	2.12
SPLIT-AUX	9034	81.66	3.89	0.57
SPLIT-CC	9190	81.69	3.92	0.12
SPLIT-%	9255	81.81	4.04	0.15
TMP-NP	9594	82.25	4.48	1.07
GAPPED-S	9741	82.28	4.51	0.17
POSS-NP	9820	83.06	5.29	0.28
SPLIT-VP	10499	85.72	7.95	1.36
BASE-NP	11660	86.04	8.27	0.73
DOMINATES-V	14097	86.91	9.14	1.42
RIGHT-REC-NP	15276	87.04	9.27	1.94

Figure 11.10: Parsing accuracies as new non-terminals are introduced

Automated state-splitting Later work from Petrov et al. (2006) automated the state-splitting process, using expectation maximization.

- Assign a random subcategory to each node.
- Learn a PCFG.
- Apply the PCFG to relabel the nodes
 - subject to constraints of original annotations:
VP3 can be relabeled as VP7, but not as an NP
- Repeat

They applied a split-merge heuristic to determine the number of subcategories for each phrase type. See slides for more details.

Conditional Random Field parsing

We can think of a PCFG parser in our usual framework of structured prediction:

$$\hat{\tau} = \arg \max_{\tau} \boldsymbol{\theta}^{\top} \mathbf{f}(\tau, \mathbf{w}). \quad (11.20)$$

In this case, the features $\mathbf{f}(\tau, \mathbf{w})$ count all the CFG productions in τ and the terminal productions to \mathbf{w} , and the weights $\boldsymbol{\theta}$ count the log-probabilities of those productions.

This suggests that we could try to learn the weights $\boldsymbol{\theta}$ discriminatively. But if we are willing to learn the weights discriminatively, we can also add additional features; we only require a feature decomposition so that $\mathbf{f}(\tau, \mathbf{w})$ decomposes across the productions in τ , so that we can still perform CKY parsing to find the best-scoring parse. For example, under such a decomposition, we could incorporate lexical features, so that we learn weights for the non-terminal production as well as for lexicalized forms,

$f1 \text{ NP}(*) \rightarrow \text{NP}(*) \text{PP}(*)$

$f2 \text{ NP}(cats) \rightarrow \text{NP}(cats) \text{PP}(*)$

$f3 \text{ NP}(*) \rightarrow \text{NP}(*) \text{PP}(claws)$

$f4 \text{ NP}(cats) \rightarrow \text{NP}(cats) \text{PP}(claws)$

Through regularization, we can find weights that strike a good balance between frequently-observed features ($f1$) and more discriminative features ($f4$).

This approach was implemented by Finkel et al. (2008) in the context of PCFG parsing with Conditional Random Fields. They used stochastic gradient descent for training, with the inside-outside algorithm (analogous to forward-backward, but for trees) to compute expected feature counts. However, the time complexity of $\mathcal{O}(M^3)$ posed serious challenges — recall that CRF sequence labeling can be trained in linear time. Finkel et al. (2008) address these issues by “prefiltering” the CKY parsing chart, identifying the productions which cannot be part of any complete parse.

Carreras et al. (2008) use the averaged perceptron to perform conditional parsing, employing an alternative feature decomposition based on tree-adjoining grammar (TAG). This yields substantially better results ($F = 90.5$), but is still less accurate than the reranked generative parser of McClosky et al. (2006).

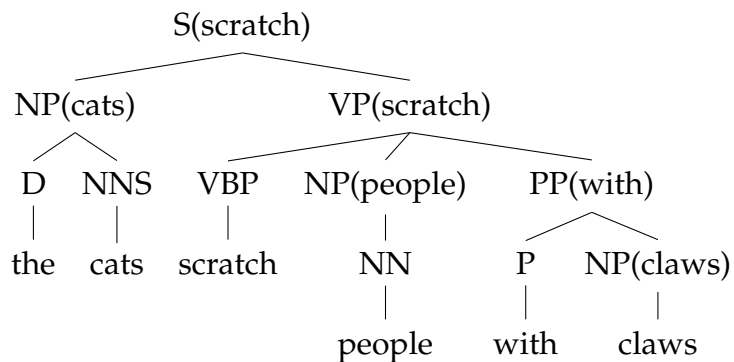
Vanilla PCFG	72%
Parent-annotations (Johnson, 1998)	80%
Lexicalized (Charniak, 1997)	86%
Lexicalized (Collins, 2003)	87%
Lexicalized, reranking, self-training (McClosky et al., 2006)	92.1%
State splitting (Petrov and Klein, 2007)	90.1%
CRF Parsing (Finkel et al., 2008)	89%
TAG Perceptron Parsing (Carreras et al., 2008)	90.5%
Compositional Vector Grammars (Socher et al., 2013)	90.4%

Table 11.5: Parsing scoreboard, circa 2013(Socher et al., 2013)

Chapter 12

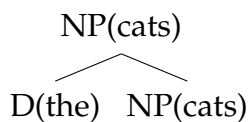
Dependency Parsing

Lexicalized parsing augments the non-terminals with **head words**.

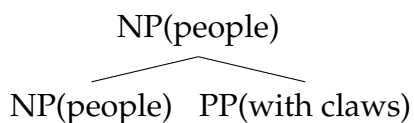


- A set of deterministic **head percolation** rules determine how heads move up the tree through each production.
- Some rules are pretty obvious:

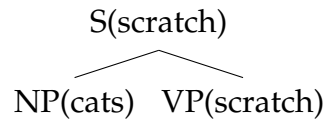
- the head of a determiner-noun constituent is the noun:



- prepositional adjuncts are not heads:

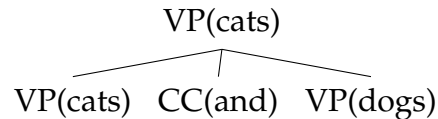


- verbal predicates are the heads of sentences:

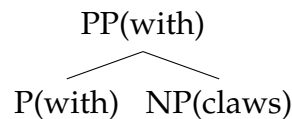


- Others are less clear-cut:

- the head of a conjunction is the left element:



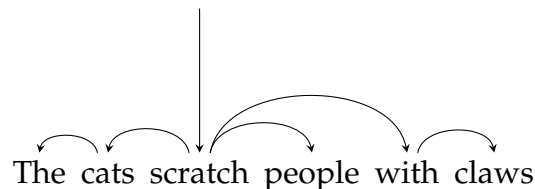
- the head of a prepositional phrase is the preposition:



(Alternatively, we can “collapse” out the preposition itself. This is the approach taken by the Stanford dependency parser.)

- We could argue about this stuff, but once we agree on a standard it can be applied deterministically to any parse tree.

A head-annotated parse tree defines a graph over the words in the sentence:

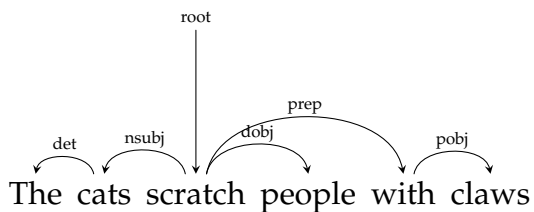


What are the properties of this graph?

- directed
- weakly connected
- every node has one incoming edge
- (therefore) no cycles
- (therefore) a tree

What is the meaning of the edges?

- A dependency edge means there is an asymmetric syntactic relationship between the head and the modifier.
- (sometimes called nucleus/satellite, or governor/dependent, or parent/child)
- Criteria for figuring out who is the head:
 - The modifier may be optional while the head is mandatory: (*red HAT*, *EAT quickly*)
 - The head sets the syntactic category of the construction: (prepositions are the heads of prepositional phrases)
 - The head determines the morphological form of the modifier: (*los LI-BROS*, *una CASA*, *these HOUSES*)
- As always, these guidelines sometimes conflict.
- Edges may be **labeled** to indicate their function:

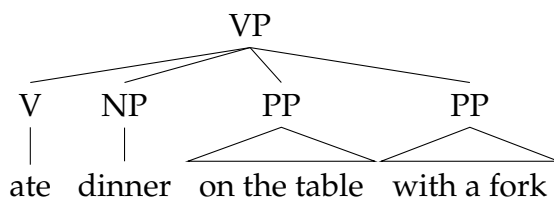


Dependency trees tell us who did what to whom.

Expressiveness

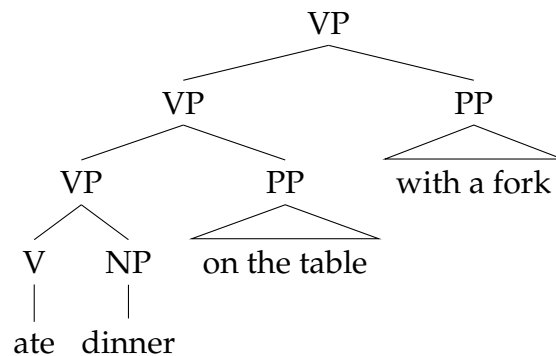
(Unlabeled) dependency trees are less expressive than CFG derivations. That means they hide some of the ambiguity in CFGs that we may not care about. Remember the different representations for PP modification?

- Flat

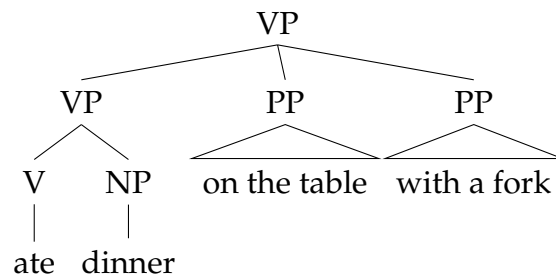


(c) Jacob Eisenstein 2014-2015. Work in progress.

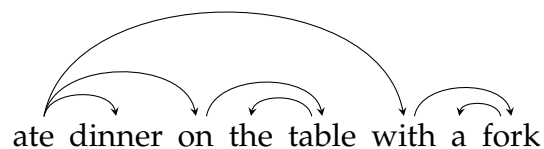
- (Chomsky) adjunction



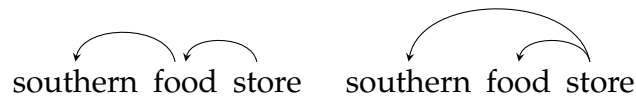
- Two-level (PTB)



These all look the same in a dependency parse:



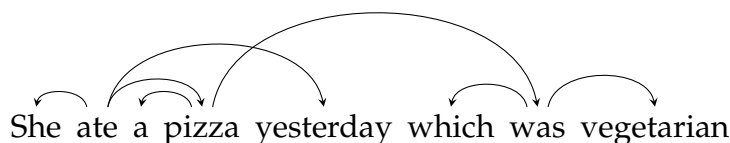
So if you didn't think there was any meaningful difference between these representations, you should be happy. But many kinds of CFG ambiguity remain in the dependency parse:



Projectivity

In projective dependency parsing, there can be no crossing edges in the dependency graph.

- More formally, for every word w_h , there must be a span $i : j$ such a word w_m is a descendant of w_h if and only if it lies within the span w_i, w_{i+1}, \dots, w_j .
- Crossing edges are rare in English:



However, they are more common in other languages, like Czech.

- We can build projective dependency banks from constituent treebanks, like PTB, by using a set of head-finding rules.
 - In this case, crossing edges are prohibited by construction.
 - What is the CFG analysis for the pizza example?
- In languages where non-projectivity is common, we must annotate dependency trees directly. An example is the Prague dependency Treebank, which contains 1.5M words of Czech.

12.1 Algorithms

Let's assume that the score for a dependency parse can be **factored** across the edges:

$$\Psi(G) = \bigotimes_i \psi(g(i) \rightarrow i, r), \quad (12.1)$$

where $g(i)$ is the parent of node i and r is the relation type. If we can assume this factorization, then we have efficient dynamic programs for dependency parsing. The remainder of this section will concern itself with unlabeled dependency parsing, dropping the label r .

Maximum spanning tree

Non-projective dependency parsing reduces to the maximum spanning tree problem (in directed graphs).

- We can build a connected graph, with edge weights equal to $\psi(i \rightarrow j)$ for all words i, j in the sentence. In the case of labeled parsing, we can build

multiple edges between each pair of words, one for each relation r , with weight $\psi(i \rightarrow j, r)$.

- To get a dependency parse, we need a tree that touches all the nodes.
- To get the **best** dependency parse, we need a tree that achieves the maximum score $\bigotimes_i \psi(g(i) \rightarrow i)$, where $g(i)$ is the head of node i .
- The Chu-Liu-Edmonds algorithm computes this in $\mathcal{O}(N^3)$. [See slides.]
- The Tarjan algorithm is $\mathcal{O}(N^2)$.

Dynamic programs

Projective dependency parsing can be performed with dynamic programming. Here is a naive recursive algorithm, for unlabeled dependency parsing.

- $c[i, j, h]$ is the score of the best tree spanning $i \rightarrow j$ with head h

$$c[i, j, h] = \left(\bigoplus_{k, h'} c[i, k, h] \otimes c[k, j, h'] \otimes \psi(h \rightarrow h') \right) \\ \oplus \left(\bigoplus_{k, h'} c[i, k, h'] \otimes c[k, j, h] \otimes \psi(h \rightarrow h') \right)$$

- The first line represents left-branching trees; the second line represents right-branching trees. Projectivity guarantees that we can find a span $[i, j]$ such that all the words $w_{i:j}$ are children of h .
- **What is the complexity?** The size of the table is $\mathcal{O}(N^3)$. To fill in each node, we must consider $\mathcal{O}(N)$ possible split points and $\mathcal{O}(N)$ possible heads h' for the satellite subtree. So the time complexity is $\mathcal{O}(N^5)$.
- The Eisner algorithm reduces this to $\mathcal{O}(N^3)$, by adapting the CKY algorithm. We keep four tables instead of 1!
 - scores of **incomplete** subtrees from i to j , headed to the left
 - scores of **incomplete** subtrees from i to j , headed to the right
 - scores of **complete** subtrees from i to j , headed to the left
 - scores of **complete** subtrees from i to j , headed to the right
- Why?

- Incomplete subtrees can subsume complete subtrees heading in the same direction, resulting in a complete subtree.
- Complete subtrees can combine if they are heading in opposite directions, resulting in an incomplete subtree.
- Our goal is to produce a complete tree from 0 to M .

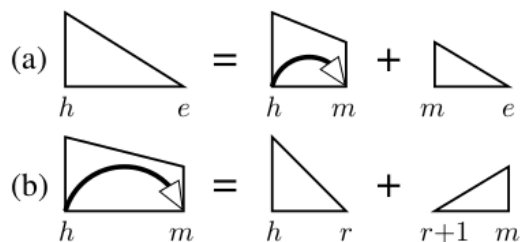


Figure 12.1: Diagram of Eisner algorithm for first-order dependency parsing (Koo and Collins, 2010)

The Eisner algorithm in action

Let's consider the unlabeled projective dependency parse for the phrase “plastic cup holders.”

Assume we have the following log potentials $\psi_{i,j} = \boldsymbol{\theta}^\top \mathbf{f}(i, j)$.

	ROOT	plastic	cup	holders
Edge weights:	ROOT	1	1	1
	plastic	$-\infty$	-1	-1
	cup	$-\infty$	2	-1
	holders	$-\infty$	0	4

Let's assume we have a 4-dimensional table C , such that $C[i, j, d, c]$ is the score of:

- the best subtree from i to j , with $i, j \in [1, M]$
- in direction $d \in \{\leftarrow, \rightarrow\}$
- with completeness $c \in \{0, 1\}$

(c) Jacob Eisenstein 2014-2015. Work in progress.

If $d = \leftarrow$ then the head of the tree is j , the right-most element; otherwise the head is the left-most element, i . If completeness $c = 1$, then the subtree is not taking any more dependents; otherwise it needs to be completed.

The dynamic program works as follows

- Build an incomplete subtree by merging two adjacent subtrees a and b , where a is right-facing and b is left-facing. The new subtree is left-facing if we add a dependency from the right edge of b to the left edge of a ; otherwise it's right-facing.
- Build a complete left-facing subtree by merging a complete left-facing subtree a with an adjacent incomplete left-facing subtree b .
- Build a complete right-facing subtree by merging an incomplete right-facing subtree a with a complete left-facing subtree b .

At each merge, we search for the best split point, which is scored by adding the scores of the subtrees. Specifically:

$$\begin{aligned}
 C[s, t, \leftarrow, 0] &= \psi_{t,s} + \max_{s \leq r < t} (c[s, r, \rightarrow, 1] + c[r + 1, t, \leftarrow, 1]) \\
 C[s, t, \rightarrow, 0] &= \psi_{s,t} + \max_{s \leq r < t} (c[s, r, \rightarrow, 1] + c[r + 1, t, \leftarrow, 1]) \\
 C[s, t, \leftarrow, 1] &= \max_{s \leq r < t} (c[s, r, \leftarrow, 1] + c[r, t, \leftarrow, 0]) \\
 C[s, t, \rightarrow, 1] &= \max_{s \leq r < t} (c[s, r, \rightarrow, 0] + c[r, t, \rightarrow, 1])
 \end{aligned}$$

We also need to keep back pointers. The score of the final parse is $C[1, n, \rightarrow, 1]$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

$$\begin{aligned}
c[1, 2, \leftarrow, 0] &= c[1, 1, \rightarrow, 1] + c[2, 2, \leftarrow, 1] + \psi(2, 1) = \psi(2, 1) = -\infty \\
c[1, 2, \rightarrow, 0] &= c[1, 1, \rightarrow, 1] + c[2, 2, \leftarrow, 1] + \psi(1, 2) = \psi(1, 2) = 1 \\
c[1, 2, \leftarrow, 1] &= c[1, 1, \leftarrow, 1] + c[1, 2, \leftarrow, 0] = c[1, 2, \leftarrow, 0] = \psi(2, 1) = -\infty \\
c[1, 2, \rightarrow, 1] &= c[1, 2, \rightarrow, 0] + c[2, 2, \rightarrow, 1] = c[1, 2, \rightarrow, 0] = \psi(1, 2) = 1
\end{aligned}$$

$$\begin{aligned}
c[2, 3, \leftarrow, 0] &= \psi(3, 2) = 2 \\
c[2, 3, \rightarrow, 0] &= \psi(2, 3) = -1 \\
c[2, 3, \leftarrow, 1] &= c[2, 3, \leftarrow, 0] = \psi(3, 2) = 2 \\
c[2, 3, \rightarrow, 1] &= \psi(2, 3) = -1
\end{aligned}$$

$$\begin{aligned}
c[3, 4, \leftarrow, 0] &= \psi(4, 3) = 4 \\
c[3, 4, \rightarrow, 0] &= \psi(3, 4) = -1 \\
c[3, 4, \leftarrow, 1] &= c[3, 4, \leftarrow, 0] = \psi(4, 3) = 4 \\
c[3, 4, \rightarrow, 1] &= \psi(3, 4) = -1
\end{aligned}$$

$$\begin{aligned}
c[1, 3, \leftarrow, 0] &= \psi(3, 1) + \max(c[1, 1, \rightarrow, 1] + c[2, 3, \leftarrow, 1], c[1, 2, \rightarrow, 1] + c[3, 3, \leftarrow, 1]) \\
&= \psi(3, 1) + \max(\psi(3, 2), \psi(1, 2)) = \psi(3, 1) + \psi(3, 2) = -\infty \\
c[1, 3, \rightarrow, 0] &= \psi(1, 3) + \max(c[2, 3, \leftarrow, 1], c[1, 2, \rightarrow, 1]) = \psi(1, 3) + \psi(3, 2) = 3 \\
c[1, 3, \leftarrow, 1] &= \max(c[1, 1, \leftarrow, 1] + c[1, 3, \leftarrow, 0], c[1, 2, \leftarrow, 1] + c[2, 3, \leftarrow, 0]) = \max(0 - \infty, -\infty + 2) = -\infty \\
c[1, 3, \rightarrow, 1] &= \max(c[1, 2, \rightarrow, 0] + c[2, 3, \rightarrow, 1], c[1, 3, \rightarrow, 0] + c[3, 3, \rightarrow, 1]) \\
&= \max(\psi(1, 2) + \psi(2, 3), \psi(\mathbf{1}, \mathbf{3}) + \psi(\mathbf{3}, \mathbf{2})) = \max(0, 3) = 3
\end{aligned}$$

$$\begin{aligned}
c[2, 4, \leftarrow, 0] &= \psi(4, 2) + \max(c[3, 4, \leftarrow, 1], c[2, 3, \rightarrow, 1]) = \psi(4, 2) + \max(\psi(\mathbf{4}, \mathbf{3}), \psi(2, 3)) = 4 \\
c[2, 4, \rightarrow, 0] &= \psi(2, 4) + \max(\psi(\mathbf{4}, \mathbf{3}), \psi(2, 3)) = 3 \\
c[2, 4, \leftarrow, 1] &= \max(c[2, 4, \leftarrow, 0], c[2, 3, \leftarrow, 1] + c[3, 4, \leftarrow, 0]) \\
&= \max(\psi(4, 2) + \psi(4, 3), \psi(\mathbf{3}, \mathbf{2}) + \psi(\mathbf{4}, \mathbf{3})) = 6 \\
c[2, 4, \rightarrow, 1] &= \max(c[2, 3, \rightarrow, 0] + c[3, 4, \rightarrow, 1], c[2, 4, \rightarrow, 0]) \\
&= \max(\psi(2, 3) + \psi(3, 4), \psi(\mathbf{2}, \mathbf{4}) + \psi(\mathbf{4}, \mathbf{3})) = 3
\end{aligned}$$

$$\begin{aligned}
c[1, 4, \leftarrow, 0] &= \psi(4, 1) + \dots = -\infty \\
c[1, 4, \rightarrow, 0] &= \psi(1, 4) + \max(c[2, 4, \leftarrow, 1], c[1, 2, \rightarrow, 1] + c[3, 4, \leftarrow, 1], c[1, 3, \rightarrow, 1]) \\
&= \psi(1, 4) + \max(\psi(\mathbf{3}, \mathbf{2}) + \psi(\mathbf{4}, \mathbf{3}), \psi(1, 2) + \psi(4, 3), \psi(1, 3) + \psi(3, 2)) = 1 + 6 = 7 \\
c[1, 4, \leftarrow, 1] &= \max(0 + -\infty, -\infty + \dots, -\infty + \dots) = -\text{infty} \\
c[1, 4, \rightarrow, 1] &= \max(c[1, 2, \rightarrow, 0] + c[2, 4, \rightarrow, 1], c[1, 3, \rightarrow, 0] + c[3, 4, \rightarrow, 1], c[1, 4, \rightarrow, 0]) \\
&= \max(\psi(1, 2) + \psi(2, 4) + \psi(4, 3), \psi(1, 3) + \psi(3, 2) + \psi(3, 4), \psi(\mathbf{1}, \mathbf{4}) + \psi(\mathbf{4}, \mathbf{3}) + \psi(\mathbf{3}, \mathbf{2})) \\
&= \max(1 + 3, 3 - 1, 7) = 7
\end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

This corresponds to the dependency parse *plastic* < (*cup* < *holders*)

Transition-based parsing

An alternative to exact global inference is transition-based parsing: making a series of local decisions. We can apply a shift-reduce algorithm, just as we considered for CFG parsing. The reduce actions are different: rather than combining elements into non-terminals, they create arcs between words, leaving the head of edge.

- Read the sentences left-to-right,
 - **shift**: push a word onto the stack
 - **right-reduce**: make a right-facing edge between the top two elements on the stack
 - **left-reduce**: make a left-facing edge between the top two elements on the stack
 - Alternatively, “arc-eager” dependency parsing distinguishes **reduce** from **arc-right** and **arc-left**, which create arcs between the top of the stack and the first element in the queue. Arc-eager parsing is arguably more cognitively plausible, because it constructs larger connected components incrementally, rather than having a deep stack with lots of disconnected elements (Abney and Johnson, 1991; Nivre, 2004).
- **Beam search** is an improvement on shift-reduce.
 - We keep a “beam” of possible hypotheses.
 - At each stage, let the k best unique hypotheses stay on the beam.
 - In this way, we are robust to greedy bad decisions.
- **Learning**
 - Identify the series of decisions required to produce the correct dependency parse.
 - Each decision in the derivation of the correct parse is a positive instance. Each other possible decision is a negative instance.
 - Huang et al. (2012) offer alternative perceptron learning rules that yield improvements when learning in the beam search setting.

A key advantage of transition-based parsing is that there is no restriction to arc-factored features; we can include any feature of the current partial parse, history of decisions, etc. It is also fast: linear time in the length of the sentence.

12.2 Higher-order dependency parsing

Arc-factored dependency parsers can only score dependency graphs as a product across their edges. Higher-order parsers (Koo and Collins, 2010) are able to consider pairs or triples of edges (Figure 12.2)

- Second-order features consider **siblings** and **grandparents**.
- Third-order features consider **grand-siblings** (siblings and grandparents together) and **tri-siblings**.

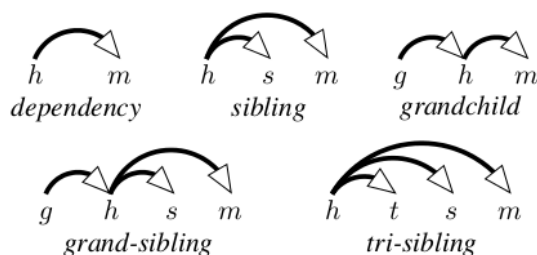


Figure 12.2: Feature templates for higher-order dependency parsing (Koo and Collins, 2010)

Why might we need higher-order dependency features? Again consider the example *cats scratch people with claws*, where the preposition *with* could attach to either *scratch* or *people*. In a lexicalized first-order arc-factored dependency parser, we would have the following feature sets for the two possible parses:

- $\langle \text{ROOT} \rightarrow \text{scratch} \rangle, \langle \text{scratch} \rightarrow \text{cats} \rangle, \langle \text{scratch} \rightarrow \text{people} \rangle, \langle \text{scratch} \rightarrow \text{with} \rangle, \langle \text{with} \rightarrow \text{claws} \rangle$
- $\langle \text{ROOT} \rightarrow \text{scratch} \rangle, \langle \text{scratch} \rightarrow \text{cats} \rangle, \langle \text{scratch} \rightarrow \text{people} \rangle, \langle \text{people} \rightarrow \text{with} \rangle, \langle \text{with} \rightarrow \text{claws} \rangle$

The only difference between the feature vectors are the features $\langle \text{scratch} \rightarrow \text{with} \rangle$ and $\langle \text{people} \rightarrow \text{with} \rangle$, but both are reasonable features, both syntactically and semantically. A first-order arc-factored dependency parsing model would therefore

struggle to find the right solution to this sentence. However, if we add grandparent features, then our feature sets include:

- $\langle \textit{scratch} \rightarrow \textit{with} \rightarrow \textit{claws} \rangle$
- $\langle \textit{people} \rightarrow \textit{with} \rightarrow \textit{claws} \rangle$,

The first feature is preferable, so a second-order dependency parser would have a better chance of correctly parsing this sentence.

Projective second-order parsing can still be performed in $\mathcal{O}(M^3)$ time (and $\mathcal{O}(M^2)$ space), using a modified version of the Eisner algorithm that includes “sibling spans.” Projective third-order parsing can be performed in $\mathcal{O}(M^4)$ time and $\mathcal{O}(M^3)$ space. Non-projective second-order dependency parsing is NP-Hard (Neuhaus and Bröcker, 1997), by reduction from the vertex cover problem. One approach is to do projective parsing first, and then post-process the projective dependency parse to add non-projective edges (Nivre and Nilsson, 2005).

12.3 Learning dependency parsers

Returning to arc-factored dependency parsing, we can easily design both generative and discriminative models:

- Generative: $\psi(m \rightarrow n) = \log P(h = m|n)$, estimated from a treebank.
- Log-linear: $\psi(m \rightarrow n) = \boldsymbol{\theta}^\top \mathbf{f}(w_m, w_n)$. Features:
 - POS tag of w_m and w_n
 - Word identity of w_m and w_n
 - Word shape (e.g., suffix, prefix)
 - Distance ($m - n$)
 - ...

Structured perceptron

In the log-linear model above, it is easy to learn the weights using structured perceptron

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}' \in \mathcal{T}(\mathbf{w})} \boldsymbol{\theta}^\top \mathbf{f}(\mathbf{w}, \mathbf{y}')$$

$$\boldsymbol{\theta}^\top = \boldsymbol{\theta}^\top + \mathbf{f}(\mathbf{w}, \mathbf{y}) - \mathbf{f}(\mathbf{w}, \hat{\mathbf{y}})$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

This is just like sequence labeling, but now $\arg \max_{y' \in \mathcal{T}(x)}$ searches over dependency trees, using either MST or the Eisner algorithm.

Conditional Random Fields (CRFs)

CRFs are just globally-normalized conditional models, and they can also be applied to any graphical model in which we can efficiently compute marginals. The prediction step is identical to the structured perceptron (using MST or the Eisner algorithm); for learning, we have a gradient that is analogous to the case in sequence labeling:

$$\frac{\partial \mathcal{L}}{\partial \theta} = \sum_m \mathbf{f}(x_{h(m)} \rightarrow w_m) - \sum_n P(n \rightarrow m | \mathbf{w}) \mathbf{f}(w_n \rightarrow w_m) \quad (12.2)$$

We require **marginal** probabilities $P(n \rightarrow m | \mathbf{w})$. These can be obtained efficiently using a variant of inside-outside (for projective trees), and the matrix-tree theorem for non-projective trees (Koo et al., 2007).

12.4 Applications

Dependency parsing is used in many real-world applications: any time you want to know about pairs of words which might not be adjacent, you can use dependency links instead of typical regular expression search patterns. For example, we may want to match strings like *delicious pastries*, *delicious French pastries*, and *the pastries are delicious*¹

- It is now possible to search Google n-grams by dependency edges.
- Muralidharan and Hearst (2013) show how dependency parsing can be used in humanities research.
- Cui et al. (2005) show how dependency parsing can improve question answering. For example:
 - Question: *What % of the nation's cheese does Wisconsin produce?*
 - Now suppose your corpus contains this sentence: *In Wisconsin, where farmers produce 28% of the nation's cheese, ...*

¹Note that the copula *is* is collapsed in many dependency parsing systems, such as the Stanford dependency parser De Marneffe and Manning (2008).

- The location of *Wisconsin* in the surface form of this string might make it a poor match for the query. However, in the dependency graph, there is an edge from *produce* to *Wisconsin* in both the question and the potential answer, raising the likelihood that this span of text is relevant to the question.
- In sentiment analysis, the polarity of a sentence can be reversed by negation, e.g. *There is no reason at all to believe the polluters will suddenly become reasonable*. By tracking the sentiment polarity through the dependency parse, we can better identify the overall polarity of the sentence (Wilson et al., 2005; Nakagawa et al., 2010).

Part III

Meaning

Chapter 13

Logical semantics

A grand ambition of natural language processing, and indeed, all of artificial intelligence, is to convert natural language into a representation that supports **semantic inferences**.¹ Jurafsky and Martin (2009) compare several alternative representations, showing parallels between representations that are apparently quite distinct. Therefore, we will focus on logical representations, in particularly **first-order logic** and **lambda calculus**.

There are several properties that we might desire in a meaning representation.

Truth conditional The “meaning” of an assertion is identical to its truth conditions.

Model-theoretic Individual elements in the meaning representation are grounded in elements in a “model” of the world, which can be thought of as a sort of database.

Inference We would like to be able to combine assertions to infer new facts about the world.

...

The goal of **semantic parsing** is to convert natural language statements to a representation that meets these criteria. For example, given a statement like *all pastries are delicious*, we would parse to the first-order logical representation:

$$\forall x : \text{PASTRIES}(x) \implies \text{DELICIOUS}(x) \quad (13.1)$$

¹This chapter is just a sketch. In class we have used the chapter from Jurafsky and Martin (2009), and a more involved “informal” reading from Levy and Manning (2009); another possible introduction is from Briscoe (2011).

There are several elements here:

- The **operator** \Rightarrow , which you may recognize from Boolean logic as representing logical implication. An operator converts one or more **truth values** into another truth value; in this case $\alpha \Rightarrow \beta$ is **true** unless α is true and β is false.
- The **predicates** PASTRIES and DELICIOUS. In a model-theoretic semantics, these can be thought of as representing set membership, where $\text{PRED}(x)$ is **true** iff the element x is in the set PRED.
- The **variable** x , which is **bound** by the **quantifier** \forall . To assess the truth value of an assertion, we will need all variables to be bound.

Another type of element that we will need are **constants**, which correspond to specific objects in the model. This enables us to analyze sentences like *Umashanthi likes jazz*, where the **denotation** $[\text{Umashanthi}] = u$, where u is some object in the model — think of a single record in a database, for example.

To come:

- Compositionality
- Lambda calculus
- Syntax-semantic interface (CFG-style or CCG?)
- Semantic types
- Ambiguity
- Determiners?
- Question answering?
- Learning semantic parsers

Chapter 14

Shallow semantics

(This is just dumped from the slides, needs lots of work. Is also a little outdated now.)

“Full” compositional semantics requires representations at least as expressive as first-order logic. Machine learning approaches have improved robustness, and recent work has driven down the requirements for manually-created resources. But coverage is still relatively limited, with best performance in narrow domains like travel and geography.

Shallow semantics comprises a set of alternative approaches, which trade the expressiveness of representations like first-order logic for shallower representations which can be parsed more robustly, with broader coverage.

14.1 Predicates and arguments

Shallow semantics focuses on predicate-argument relations. For example, the sentence *Regina trusts Ben* can be interpreted as `trusts(REGINA, BEN)`, where `trusts` is a predicate and `REGINA` and `BEN` are its arguments. This is exactly the sort of relation that we saw in first-order logical (FOL) semantics too, but in shallow semantics we will typically work without variables and quantification. (However, there have been recent explorations intermediate representations between FOL and shallow predicate-argument relations, as described in Section 14.5.)

To see how shallow semantics can represent meaning, consider these four sentences (borrowed from the slides of a tutorial by Kristina Toutanova and Scott Yih).

- Yesterday, Kristina hit Scott with a baseball

- Scott was hit by Kristina yesterday with a baseball
- Yesterday, Scott was hit with a baseball by Kristina
- Kristina hit Scott with a baseball yesterday

We don't need first-order logic to realize that these sentences are semantically identical. Shallow semantics will suffice: the *roles* in each sentence are filled by the same text.

- **Hitter**: Kristina
- **Person hit**: Scott
- **Instrument of hitting**: with a baseball
- **Time of hitting**: yesterday

Deep roles The event semantics representation for the sentence *Scott was hit by Kristina yesterday* (and all of the other examples) is:

$$\begin{aligned} \exists e, x, y \text{ Hitting}(e) \wedge \text{Hitter}(e, \text{Kristina}) \wedge \text{PersonHit}(e, \text{Scott}) \\ \wedge \text{TimeOfHitting}(e, \text{Yesterday}) \end{aligned}$$

In this example, *Hitter*, *PersonHit*, and *TimeOfHitting* are roles. We use these specific roles because of the **predicate verb** *hit*. Roles that relate to a specific predicate are called “deep roles.”

Thematic roles Without knowing more about deep roles like *Hitter*, we cannot do much inference. But building classifiers for every role of every predicate would be a lot of work, and we would struggle to get enough training data to accomplish this. Is there a shortcut?

Consider the example *Scott was paid by Kristina yesterday*. Arguably, the role-fillers *Scott*, *Kristina* and *yesterday* have similar thematic functions as in the earlier sentence about baseballs. **Thematic roles** attempt to capture the similarity between *Payer* and *Hitter*, and between *PersonHit* and *PersonPaid*.¹

Here is a table of some typical thematic roles.²

¹Thematic roles date to Panini (7th-4th century BCE!). The modern formulation is due to Fillmore (1968) and Gruber (1965).

²These examples might have been borrowed from K&S's slides, try to track this down.

AGENT	The volitional causer <i>The waiter spilled the soup</i>
EXPERIENCER	The experiencer <i>The soup gave all three of us a headache.</i>
FORCE	The non-volitional causer <i>The wind blew my soup off the table.</i>
THEME	The participant most directly affected <i>The wind blew my my soup off the table.</i>
RESULT	The end product <i>The cook has prepared a cold duck soup.</i>
CONTENT	The proposition or content of a propositional event <i>The waiter assured me that the soup is vegetarian.</i>
INSTRUMENT	An instrument used in an event <i>It's hard to eat soup with chopsticks.</i>
BENEFICIARY	The beneficiary <i>The waiter brought me some soup.</i>
SOURCE	The origin of the object of a transfer event <i>The stack of canned soup comes from Pittsburgh.</i>
GOAL	The destination of the object of a transfer event <i>He brought the bowl of soup to our table.</i>

Case frames Different verbs take different thematic roles as arguments. The possible arguments for a verb is the **case frame** or **thematic grid**. For example, for *break*:

- AGENT: Subject, THEME: Object
John broke the window.
- AGENT: Subject, THEME: Object, INSTRUMENT: PP (with)
John broke the window with a rock.

- INSTRUMENT: Subject, THEME: Object
The rock broke the window.
- THEME: Subject
The window broke.

When two verbs have similar case frames, this is a clue that they might be semantically related: (e.g., *break*, *shatter*, *smash*).

Many verbs permit multiple orderings of the same arguments. These are known as **diathesis alternations**. For example, *give* permits the dative alternation,

[*agent* Doris] **gave** [*goal* Cary] [*theme* the book].

which could also be written,

[*agent* Doris] **gave** [*theme* the book] [*goal* to Cary]

Again, similar alternation patterns suggest semantic similarity. For example, verbs that display the dative alternation include some broad classes:

- “verbs of future having” (advance, allocate, offer, owe)
- “verbs of sending” (forward, hand, mail)
- “verbs of throwing” (kick, pass, throw)

The purpose of thematic roles is to abstract above verb-specific roles. But it is usually possible to construct examples in which thematic roles are insufficiently specific.

- *Intermediary instruments* can act as subjects:
 1. *The cook opened the jar with the new gadget.*
 2. *The new gadget opened the jar.*
- *Enabling instruments* cannot:
 1. *Shelly ate the pizza with the fork.*
 2. **The fork ate the pizza.*

Thematic roles are bundles of semantic properties, but it’s not clear how many properties are necessary. For example, AGENTS are usually animate, volitional, sentient, causal, but any of these properties may be missing occasionally.

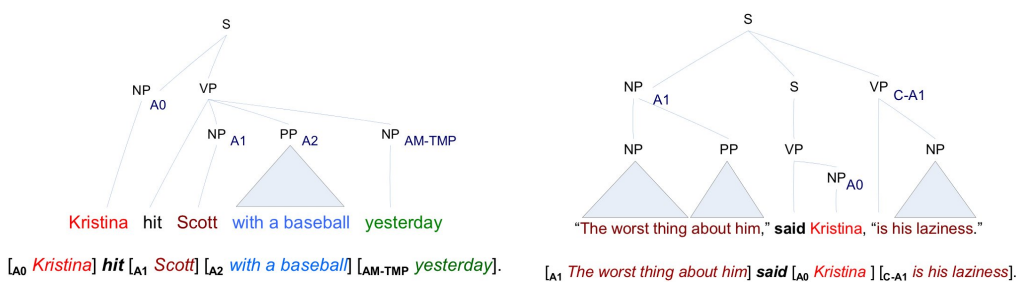


Figure 14.1: Examples of PropBank-style annotations, borrowed from the slides of Toutanova and Yih

14.2 Resources for shallow semantics

In the Proposition Bank (**PropBank**), roles are verb-specific, with some sharing Palmer et al. (2005).

- Arg0: proto-agent (has agent-like properties)
- Arg1: proto-patient (has patient-like properties)
- Arg2... ArgN: verb-specific
- 13 universal adjunct-like arguments: temporal, manner, location, cause, negation, ...

PropBank contains two main resources:

- a set of labeled sentences, built on the Penn TreeBank
- a set of “Frame Files” describing each verbal predicate

<http://verbs.colorado.edu/propbank/framesets-english/scratch-v.html>

Some example PropBank-style annotations are shown in Figure 14.1.

The PropBank corpus was last released on March 4, 2005. The details are:

- Verb Lexicon: 3,324 frame files
- Annotation: 113,000 propositions

PropBank has been used as the standard dataset for shared tasks on semantic role labeling (SRL). There are related corpora:

- Chinese PropBank <http://www.cis.upenn.edu/~chinese/cpb/>

(c) Jacob Eisenstein 2014-2015. Work in progress.

FRAMENET ANNOTATION:

[_{Buyer} Chuck] *bought* [_{Goods} a car] [_{Seller} from Jerry] [_{Payment} for \$1000].

[_{Seller} Jerry] *sold* [_{Goods} a car] [_{Buyer} to Chuck] [_{Payment} for \$1000].

PROPBANK ANNOTATION:

[_{Arg0} Chuck] *bought* [_{Arg1} a car] [_{Arg2} from Jerry] [_{Arg3} for \$1000].

[_{Arg0} Jerry] *sold* [_{Arg1} a car] [_{Arg2} to Chuck] [_{Arg3} for \$1000].

Figure 14.2: A comparison of framenet and propbank, from Toutanova and Yih
[**todo: I think**]

- NomBank: structure of noun phrases, e.g. [_{A0} Her] [_{REL} gift] of [_{A1} a book]
[_{A2} to John]

14.3 FrameNet

The key idea of FrameNet is to group related verbs (and nouns) into *frames*

- [_{A1} The price of bananas] increased [_{A2} 5%].
- [_{A1} The price of bananas] rose [_{A2} 5%].
- There has been a [_{A2} 5%] rise [_{A1} in the price of bananas].

The first two sentences involve different verbs; the second sentence conveys same semantics with a noun. Nonetheless, the meaning is the same. FrameNet captures this.

The relationship between Framenet and PropBank annotation is shown in Figure 14.2. A frame defines a set of *lexical units* and a set of *frame elements*:

(c) Jacob Eisenstein 2014-2015. Work in progress.

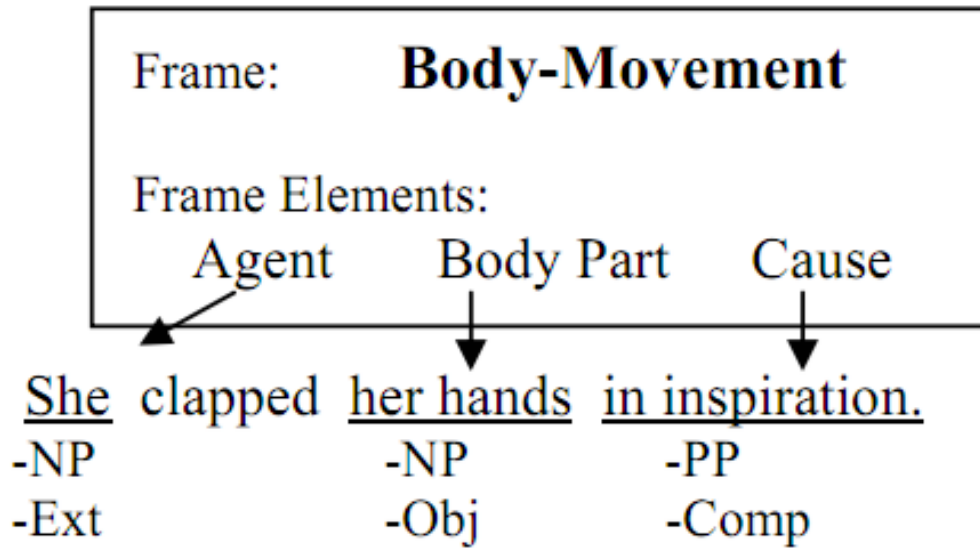
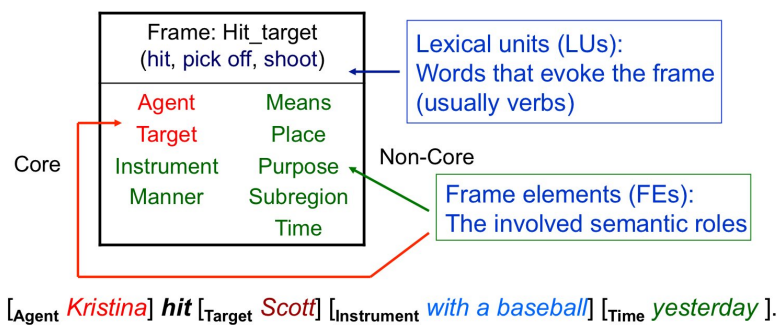


Figure 14.3: FrameNet annotation, figure from Fleischman et al, 2003



Lexical units for **Change-position-on-a-scale** frame:

VERBS:	dwindle	move	soar	escalation	shift
advance	edge	mushroom	swell	explosion	tumble
climb	explode	plummet	swing	fall	
decline	fall	reach	triple	fluctuation	ADVERBS:
decrease	fluctuate	rise	tumble	gain	increasingly
diminish	gain	rocket		growth	
dip	grow	shift	NOUNS:	hike	
double	increase	skyrocket	decline	increase	
drop	jump	slide	decrease	rise	

(c) Jacob Eisenstein 2014-2015. Work in progress.

Frame elements for **Change-position-on-a-scale** frame:

Core Roles	
ATTRIBUTE	The ATTRIBUTE is a scalar property that the ITEM possesses.
DIFFERENCE	The distance by which an ITEM changes its position on the scale.
FINAL_STATE	A description that presents the ITEM's state after the change in the ATTRIBUTE's value as an independent predication.
FINAL_VALUE	The position on the scale where the ITEM ends up.
INITIAL_STATE	A description that presents the ITEM's state before the change in the ATTRIBUTE's value as an independent predication.
INITIAL_VALUE	The initial position on the scale from which the ITEM moves away.
ITEM	The entity that has a position on the scale.
VALUE_RANGE	A portion of the scale, typically identified by its end points, along which the values of the ATTRIBUTE fluctuate.
Some Non-Core Roles	
DURATION	The length of time over which the change takes place.
SPEED	The rate of change of the VALUE.
GROUP	The GROUP in which an ITEM changes the value of an ATTRIBUTE in a specified way.

The FrameNet corpus is publicly available online: <https://framenet.icsi.berkeley.edu/fndrupal/about>. As of October 2013, they had annotated

- 1,164 semantic frames
- 12,713 lexical units
- 196,000 manually annotated sentences
- still ongoing...

Unlike PropBank, Framenet is not based on TreeBank parses, and example sentences are chosen by hand.

14.4 Semantic Role Labeling

Semantic role labeling (SRL) is the task of assigning semantic labels to spans of text. Labels describe the role of the phrase with respect to the *predicate verb*. In practice, this usually means PropBank labels, e.g. Arg0, Arg1, etc.

[**todo: example slides, continuation and reference arguments**]

(c) Jacob Eisenstein 2014-2015. Work in progress.

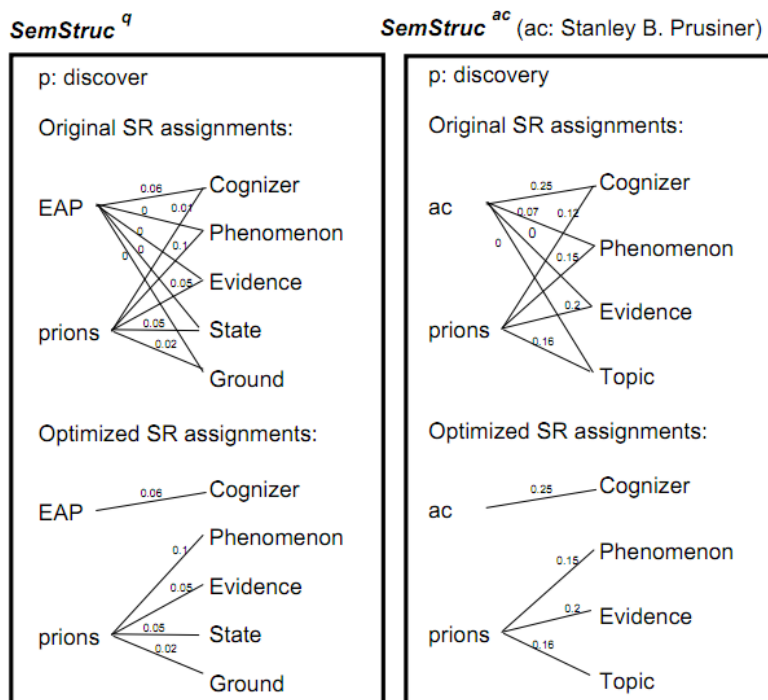


Figure 14.4: Using semantic role labeling to align questions and answers

Applications of SRL Why might we want to do this? One application is to automatic question answering systems like IBM Watson. Consider the example question, *Who discovered prions?*. Somewhere in our database, we have the statement *1997: Stanley B. Prusiner, United States, discovery of prions....* How can we link them up? Shen and Lapata (2007) use semantic roles to align questions against the content of factual sentences, as shown in Figure 14.4.

14.5 Abstract Meaning Representation

[todo: use this as an excuse to add something about dual decomposition?]

Chapter 15

Distributional semantics

A recurring theme in this course is that the mapping from words to meaning is complex.

- **Word sense disambiguation:** multiple meanings for the same form (e.g., *bank*)
- **Morphological analysis:** shared semantic basis among multiple forms (e.g., *speak, spoke, speaking*)
- **Synonymy:** in English we have lots of synonyms and near neighbors, as English combines influence from lots of other languages (French, Latin, German, etc)
- Both **compositional** and **frame** semantics assume hand-crafted resources that map from words to predicates.

How do we do semantic analysis of words that we've never seen before?

15.1 The distributional hypothesis

Here's a word you may not know: *tezgüino*. If we encounter this word, what can we do? It seems like a big problem for any NLP system, from POS tagging to semantic analysis.

Suppose we see that *tezgüino* is used in the following contexts:

1. A bottle of _____ is on the table.
2. Everybody likes _____.

3. Don't have _____ before you drive.

4. We make _____ out of corn.

What other words fit into these contexts? How about: *loud*, *motor oil*, *tortillas*, *choices*, *wine*?

We can create a vector for each word, based on whether it can be used in each context.

	C1	C2	C3	C4	...
<i>tezgüino</i>	1	1	1	1	
<i>loud</i>	0	0	0	0	
<i>motor oil</i>	1	0	0	1	
<i>tortillas</i>	0	1	0	1	
<i>choices</i>	0	1	0	0	
<i>wine</i>	1	1	1	1	

- Based on these vectors, we see:
 - *wine* is very similar to *tezgüino*
 - *motor oil* and *tortillas* are fairly similar to *tezgüino*
 - *loud* is quite different.
- The vectors describe the **distributional** properties of each word.
- Does vector similarity imply semantic similarity? This is the **distributional hypothesis**. “You shall know a word by the company it keeps.” (Firth 1957)
- It is also known as a **vector-space model**, since each word's meaning is captured by a vector.

Vector-space models and distributional semantics are relevant to a wide range of NLP applications.

- **Query expansion:** search for *bike*, match *bicycle*
- **Semi-supervised learning:** use large unlabeled datasets to acquire features which are useful in supervised learning
- **Lexicon and thesaurus induction:** automatically expand hand-crafted lexical resources, or induce them from raw text

Here are some of the practical questions that we encounter when working with vector space representations of distributional semantics:

(c) Jacob Eisenstein 2014-2015. Work in progress.

- What kinds of context should we consider? (see slides)
- How do measure similarity?
- How do we properly weigh frequent versus infrequent events?

15.2 Local context

The Brown et al (1992) clustering algorithm is over 20 years old and is still widely used in NLP!

- Context is just the immediately adjacent words.
- A generative probability model:
 - Assume each word w_i has a class c_i
 - Assume a generative model $\log p(w) = \sum_i \log p(w_i|c_i) + \log p(c_i|c_{i-1})$
(What does this remind you of?)
- Hierarchical clustering algorithm:
 - Start with every word in its own cluster
 - Until tired,
 - * Choose two clusters c_i and c_j such that merging them will give the maximum improvement in $\log p(w)$
 - * Equivalently, merge the clusters with the greatest mutual information.
 - The merge path of a word describes its semantics.

Model specifics

- \mathcal{V} is the set of all words
- N number of observed word tokens
- $n(w)$ is the number of times we see word $w \in \mathcal{V}$
- $n(w, v)$ is the number of times w precedes v
- Let $C \rightarrow \{1, 2, \dots, k\}$ define a partition of words into k classes

$$\begin{aligned} p(w_1, w_2, \dots, w_T; C) &= \prod_i p(w_i | C(w_i)) p(C(w_i) | C(w_{i-1})) \\ \log p(w_1, w_2, \dots, w_T; C) &= \sum_i \log p(w_i | C(w_i)) p(C(w_i) | C(w_{i-1})) \end{aligned}$$

This is kind of like an HMM, but each word can only be produced by a single cluster.

Let's define the "quality" of a clustering as the average log-likelihood:

$$\begin{aligned} J(C) &= \frac{1}{N} \sum_i \log (p(w_i | C(w_i)) p(C(w_i) | C(w_{i-1}))) \\ &= \sum_{w, w'} \frac{n(w, w')}{N} \log (p(w' | C(w')) p(C(w') | C(w'))) && \text{sum over word types instead} \\ &= \sum_{w, w'} \frac{n(w, w')}{N} \log \left(\frac{n(w')}{n(C(w'))} \frac{n(C(w), C(w'))}{n(C(w))} \right) && \text{definition of probabilities} \\ &= \sum_{w, w'} \frac{n(w, w')}{N} \log \left(\frac{n(w')}{1} \frac{n(C(w), C(w'))}{n(C(w)) n(C(w'))} \frac{N}{N} \right) && \text{re-arrange, multiply by one} \\ &= \sum_{w, w'} \frac{n(w, w')}{N} \log \left(\frac{n(w')}{N} \times \frac{n(C(w), C(w')) \times N}{n(C(w)) n(C(w'))} \right) && \text{re-arrange terms} \\ &= \sum_{w, w'} \frac{n(w, w')}{N} \log \frac{n(w')}{N} + \frac{n(w, w')}{N} \log \left(\frac{n(C(w), C(w')) \times N}{n(C(w)) n(C(w'))} \right) && \text{distributive law} \\ &= \sum_{w'} \frac{n(w')}{N} \log \frac{n(w')}{N} + \sum_{c, c'} \frac{n(c, c')}{N} \log \left(\frac{n(c, c') \times N}{n(c) n(c')} \right) && \text{sum across classes} \\ &= \sum_{w'} p(w') \log p(w') + \sum_{c, c'} p(c, c') \log \frac{p(c, c')}{p(c) p(c')} && \text{multiply by } \frac{N^{-2}}{N^{-2}} \text{ inside log} \\ &= -H(W) + I(C) \end{aligned}$$

In the last step, we use the following definitions from information theory:

Entropy The entropy of a discrete random variable is the expected negative log-likelihood,

$$H(X) = -E[\log P(X)] = -\sum_x P(X = x) \log P(X = x). \quad (15.1)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

For example, for a fair coin we have $H(X) = \frac{1}{2} \log \frac{1}{2} + \frac{1}{2} \log \frac{1}{2} = -\log 2$; for a (virtually) certain outcome, we have $H(x) = 1 \times \log 1 + 0 \times \log 0 = 0$. We have already seen entropy in a few other contexts.

Mutual information The information shared by two random variables is the mutual information,

$$I(X; Y) = \sum_{y \in Y} \sum_{x \in X} p_{X,Y}(x, y) \log \left(\frac{p_{X,Y}(x, y)}{p_X(x)p_Y(y)} \right). \quad (15.2)$$

For example, if X and Y are independent, then $p_{X,Y}(x, y) = p_X(x)p_Y(y)$, so the mutual information is $\log 1 = 0$. In

By $I(C)$, we are using a shorthand for the mutual information of random variables C_i and C_{i-1} — the cluster memberships of word i and $i - 1$, respectively. So we have

$$I(C) = \sum_{C_i=c, C_{i-1}=c'} \frac{P(C_i=c, C_{i-1}=c')}{P(C_i=c)P(C_{i-1}=c')} \quad (15.3)$$

The entropy $H(W)$ does not depend on the clustering, so this term is constant; choosing a clustering with maximum mutual information $I(C)$ is equivalent to maximizing the log-likelihood. Now let's see how to do that efficiently.

$V \log V$ approximate algorithm

- Take m most frequent words, put each in its own cluster c_1, c_2, \dots, c_m .
- For $i = (m + 1) : |\mathcal{V}|$
 - Create a new cluster for the c_{m+1} for word i (ordered by frequency).
 - Choose two clusters c and c' to merge, minimizing the decrease in $I(C)$. This requires $\mathcal{O}(m^2)$ operations.
- Carry out $(m - 1)$ final merges, to build full hierarchy

Cost: $\mathcal{O}(|\mathcal{V}|m^2 + n)$, plus time to sort words, $\mathcal{O}(|\mathcal{V}| \log |\mathcal{V}|)$.

15.3 Syntactic context

Local context is contingent on syntactic decisions that may have little to do with semantics:

- I gave Tim the ball.
- I gave the ball to Tim.

Using the syntactic structure of the sentence might give us a more meaningful context, yielding better clusters.

- Pereira et al (1993) cluster nouns based on the verbs for which they are the direct object.
 - The context vector for each noun is **the count of occurrences as a direct object of each verb**.
 - As with Brown clustering, a class-based probability model:

$$\begin{aligned}\hat{p}(n, v) &= \sum_{c \in \mathcal{C}} p(c, n) p(v | c) \\ &= \sum_{c \in \mathcal{C}} p(c) p(n | c) p(v | c)\end{aligned}$$

where n is the noun, v is the verb, and c is the class

- Objective: find the maximum likelihood cluster centroids.
- Dekang Lin (1997) extends this to all words, using incoming dependency edges (see slide)
 - For any pair of words i and j and relation r , we can compute:

$$p(i, j | r) = \frac{c(i, j, r)}{\sum_{i', j'} c(i', j', r)}, \quad p(i | r) = \sum_j p(i, j | r)$$

- Let $T(i)$ be the set of pairs $\langle j, r \rangle$ such that $p(i, j | r) > p(i | r)p(j | r)$
 - * $T(i)$ contains words j that are especially likely to be joined with word i in relation r .
 - * Note the connection to pointwise mutual information.
- Similarity between u and v is defined through $T(u)$ and $T(v)$.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- * Lin considers several similarity measures for $T(u)$ and $T(v)$.
- * Many of these are used widely, and are worth knowing:
 - Cosine similarity: $\frac{|T(u) \cap T(v)|}{\sqrt{|T(u)| |T(v)|}}$
 - Dice similarity: $\frac{2 \times |T(u) \cap T(v)|}{|T(u)| + |T(v)|}$
 - Jaccard similarity: $\frac{|T(u) \cap T(v)|}{|T(u)| + |T(v)| - |T(u) \cap T(v)|}$
- * Lin's metric is more complex:

$$\frac{\sum_{\langle r, w \rangle \in T(u) \cup T(v)} I(u, r, w) + I(v, r, w)}{\sum_{\langle r, w \rangle \in T(u)} I(u, r, w) + \sum_{\langle r, w \rangle \in T(v)} I(v, r, w)}$$

where $I(u, r, w)$ is the mutual information between u and w , conditioned on r .

– See slides for results.

15.4 Latent semantic analysis

(See slides)

Thus far, we have considered context vectors that are large and sparse. We can arrange these vectors into a matrix $\mathbf{X} \in \mathbb{R}^{V \times N}$, where rows correspond to words and columns correspond to contexts. However, for rare words i and j , we might have $\mathbf{x}_i^\top \mathbf{x}_j = 0$. So we'd like to have a more robust representation.

We can obtain this by factoring $\mathbf{X} \approx \mathbf{U}_K \mathbf{S}_K \mathbf{V}_K^\top$, where

$$\mathbf{U}_K \in \mathbb{R}^{V \times K}, \quad \mathbf{U}_K \mathbf{U}_K^\top = \mathbb{I} \quad (15.4)$$

$$\mathbf{S}_K \in \mathbb{R}^{K \times K}, \quad \mathbf{S}_K \text{ is diagonal} \quad (15.5)$$

$$\mathbf{V}_K \in \mathbb{R}^{D \times K}, \quad \mathbf{V}_K \mathbf{V}_K^\top = \mathbb{I} \quad (15.6)$$

Here K is a parameter that determines the fidelity of the factorization; if $K = \min(V, N)$, then $\mathbf{X} = \mathbf{U}_K \mathbf{S}_K \mathbf{V}_K^\top$. Otherwise, we have

$$\mathbf{U}_K, \mathbf{S}_K, \mathbf{V}_K = \arg \min_{\mathbf{U}, \mathbf{S}, \mathbf{V}} \|\mathbf{X} - \mathbf{U}_K \mathbf{S}_K \mathbf{V}_K^\top\|_F, \quad (15.7)$$

meaning that $\mathbf{U}_K, \mathbf{S}_K, \mathbf{V}_K$ give the rank- K matrix $\tilde{\mathbf{X}}$ that minimizes the Frobenius norm, $\sqrt{\sum_{i,j} (x_{i,j} - \tilde{x}_{i,j})^2}$.

This factorization is called Singular Value Decomposition, and is closely related to eigenvalue decomposition of the matrices $\mathbf{X}\mathbf{X}^\top$ and $\mathbf{X}^\top\mathbf{X}$. In general,

(c) Jacob Eisenstein 2014-2015. Work in progress.

the complexity of SVD is $\min(\mathcal{O}(D^2V), \mathcal{O}(V^2N))$. The standard library LAPACK (Linear Algebra PACKage) includes an iterative optimization solution for SVD, and (I think) this what is called by Matlab and Numpy.

However, for large sparse matrices it is often more efficient to take a stochastic gradient approach. Each word-context observation $\langle w, c \rangle$ gives a gradient on \mathbf{u}_w , \mathbf{v}_c , and \mathbf{S} , so we can take a gradient step. This is part of the algorithm that was used to win the Netflix challenge for predicting movie recommendation — in that case, the matrix includes raters and movies (Koren et al., 2009).

Return to NLP applications, the slides provide a nice example from (Deerwester et al., 1990), from titles of computer science research papers. In the example, the context-vector representations of the terms *user* and *human* have negative correlations, yet their distributional representations have high correlation, which is appropriate since these terms have roughly the same meaning in this dataset.

15.5 Word vectors and neural word embeddings

Discriminatively-trained word embeddings very hot area in NLP. The idea is to replace factorization approaches with discriminative training, where the task may be to predict the word given the context, or the context given the word.

Suppose we have the word w and the context c , and we define

$$u_{\theta}(w, c) = \exp(\mathbf{a}_w^{\top} \mathbf{b}_c) \quad (15.8)$$

$$(15.9)$$

with $\mathbf{a}_w \in \mathbb{R}^K$ and $\mathbf{b}_c \in \mathbb{R}^K$. The vector \mathbf{a}_w is then an **embedding** of the word w , representing its properties. We are usually less interested in the context vector \mathbf{b} ; the context can include surrounding words, and the vector \mathbf{b}_c is often formed as a sum of context embeddings for each word in a window around the current word. Mikolov et al. (2013a) draw the size of this context as a random number r .

The popular `word2vec` software¹ uses these ideas in two different types of models:

Skipgram model In the skip-gram model (Mikolov et al., 2013a), we try to maximize the log-probability of the context,

¹<https://code.google.com/p/word2vec/>

$$J = \frac{1}{T} \sum_t \sum_{-c \leq j \leq c, j \neq 0} \log p(w_{t+j} | w_j) \quad (15.10)$$

$$p(w_{t+j} | w_j) = \frac{u_\theta(w_{t+j}, w_j)}{\sum_{w'} u_\theta(w', w_j)} \quad (15.11)$$

$$= \frac{u_\theta(w_{t+j}, w_j)}{Z(w_j)} \quad (15.12)$$

This model is considered to be slower to train, but better for rare words.

CBOW The continuous bag-of-words (CBOW) (Mikolov et al., 2013b,c) is more like a language model, since we predict the probability of words given context.

$$J = \frac{1}{T} \sum_t \log p(w_t | c) \quad (15.13)$$

$$= \frac{1}{T} \sum_t \log u_\theta(w_t, c) - \log Z(c) \quad (15.14)$$

$$u_\theta(w_t, c) = \exp \left(\sum_{-c \leq j \leq c, j \neq 0} \mathbf{a}_{w_t}^\top \mathbf{b}_{w_{t+j}} \right) \quad (15.15)$$

The CBOW model is faster to train (Mikolov et al., 2013a). One efficiency improvement is build a Huffman tree over the vocabulary, so that we can compute a hierarchical version of the softmax function with time complexity $\mathcal{O}(\log V)$ rather than $\mathcal{O}(V)$. Mikolov et al. (2013a) report two-fold speedups with this approach.

These models are simplified versions of previous work on recurrent neural network language models (RNNLMs) and the log-bilinear language model (Mnih and Hinton, 2008).

Estimating word embeddings

Training these models can be challenging, because they both probabilities that need to be normalized over the entire vocabulary. This implies a training time

(c) Jacob Eisenstein 2014-2015. Work in progress.

complexity of $\mathcal{O}(VK)$ for each instance. Since these models are often trained on hundreds of billions of words, with $V \sim 10^6$ and $K \sim 10^3$, this cost is too high. Estimation techniques eliminate the factor V by making approximations.

One such approximation is negative sampling, which is a heuristic variant of noise-contrastive estimation (Gutmann and Hyvärinen, 2012).

We introduce an auxiliary variable D , where

$$D = \begin{cases} 1, & w \text{ is drawn from the empirical distribution } \hat{p}(w | c) \\ 0, & w \text{ is drawn from the noise distribution } q(w) \end{cases} \quad (15.16)$$

Now we will optimize the objective

$$\sum_{(w,c) \in \mathcal{D}} \log P(D = 1 | c, w) + \sum_{i=1, w' \sim q}^k \log P(D = 0, | c, w'), \quad (15.17)$$

setting

$$P(D = 1 | c, w) = \frac{u_\theta(w, c)}{u_\theta(w, c) + k \times q(w)} \quad (15.18)$$

$$P(D = 0 | c, w) = 1 - P(D = 1 | c, w) \quad (15.19)$$

$$= \frac{k \times q(w)}{u_\theta(w, c) + k \times q(w)}, \quad (15.20)$$

where k is the number of noise samples. Note that we have dropped the normalization $\sum_{w'} u_\theta(w', c)$. This approximation is based on the idea of noise-contrastive estimation, where the normalization term is set to be a parameter z_c . Here we approximate further, assuming $z_c = 1$. This would be trouble if we were trying to directly maximize $\log p(w | c)$, but this is where the auxiliary variable formulation helps us out: if we set θ such that $\sum_{w'} u_\theta(w' | c) \gg 1$, we will get a very low probability for $P(D = 0)$.

We can further simplify by setting $k = 1$ and $q(w)$ to a uniform distribution, arriving at

$$P(D = 1 | c, w) = \frac{u_\theta(w, c)}{u_\theta(w, c) + 1} \quad (15.21)$$

$$P(D = 0 | c, w) = \frac{1}{u_\theta(w, c) + 1} \quad (15.22)$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

The derivative with respect to a is obtained from the objective

$$L = \sum_t \log p(D = 1 \mid c_t, w_t) + \log p(D = 0 \mid c, w') \quad (15.23)$$

$$= \sum_t \log u_\theta(w_t, c_t) - \log(1 + u_\theta(w_t, c_t)) - \log(1 + u_\theta(w', c_t)) \quad (15.24)$$

$$\frac{\partial L}{\partial \mathbf{a}_i} = \sum_{t:w_t=i} \mathbf{b}_{c_t} - \frac{1}{1 + u_\theta(w_t, c_t)} \frac{\partial u_\theta(i, c_t)}{\partial \mathbf{a}_i} + \sum_t \frac{q(i)}{1 + u_\theta(i, c_t)} \frac{\partial u_\theta(i, c_t)}{\partial \mathbf{a}_i} \quad (15.25)$$

$$= \sum_{t:w_t=i} \mathbf{b}_{c_t} - P(D = 1 \mid w_t = i, c_t) \mathbf{b}_{c_t} - \sum_t q(i) P(D = 0 \mid i, c_t) \mathbf{b}_{c_t} \quad (15.26)$$

$$= \sum_t (\delta(w_t = i) - q(i)) P(D = 0 \mid w_t = i, c_t) \mathbf{b}_{c_t}. \quad (15.27)$$

The gradient with respect to \mathbf{b} is similar. In practice, we simply sample w' at each instance and compute the update with respect to \mathbf{a}_{w_t} and $\mathbf{a}_{w'}$. In practice, AdaGrad performs well for this optimization.

Connections to matrix factorization

Recent work has drawn connections between this procedure for training the skip-gram model and weighted matrix factorization approaches (Pennington et al., 2014; Levy and Goldberg, 2014). For example, Levy and Goldberg (2014) show that skip-gram with negative sampling is equivalent to factorizing a matrix M , where

$$M_{i,j} = PMI(w_i, c_j) - \log k \quad (15.28)$$

$$PMI(w_i, c_j) = \log \left(\frac{n(w = i, c = j)}{n(w = i) n(c = j)} \frac{|D|}{n(w = i) n(c = j)} \right), \quad (15.29)$$

where k is a constant offset and PMI is the well known pointwise mutual information statistic, usually written as

$$PMI(x, y) = \frac{p(x, y)}{p(x)p(y)}. \quad (15.30)$$

To see the connection, divide both the numerator and denominator by $|D|^2$.

Matrix factorization approaches are computationally advantageous because they can be applied to matrices of word co-occurrence statistics, rather than requiring streaming through an entire dataset.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Chapter 16

Discourse

16.1 Discourse relations in the Penn Discourse Treebank

- introduce discourse relations
- PDTB annotation framework in D-LTAG
- PDTB parsing

16.2 Rhetorical Structure Theory

- Higher-level discourse structure
- Shift-reduce parsing
- Applications to summarization

16.3 Centering

- Pronouns, forms of reference
- Smooth/rough transitions
- Entity grid implementation

16.4 Lexical cohesion and text segmentation

16.5 Dialogue

Minimal discussion of speech acts etc.

Chapter 17

Anaphora and Coreference Resolution

[todo: maybe blow this up and put anaphora in the discourse chapter, and coreference in the IE chapter?]

Pronouns are one of the most noticeable forms of linguistic ambiguity. A Google search for “ambiguous pronoun” reveals dozens of pages warning you to avoid ambiguity. But as we have seen, people resolve all but the most egregious linguistic ambiguities intuitively, below the level of conscious thought.

Moreover, reference ambiguities need not apply only to pronouns. Consider the following text:

Apple Inc Chief Executive Tim Cook has jetted into China for talks with government officials as he_1 seeks to clear up a pile of problems in $[[the\ firm's]_2\ biggest\ growth\ market]_3$.

- Who is referred to by he_1 ?
- What entity is referred to by $the\ firm_2$?
- What is Apple’s biggest growth market?

You probably answered these questions by making some commonsense assumptions. Tim Cook is the only individual mentioned, so the personal pronoun *he* probably refers to him; Apple is the only firm mentioned, so *the firm* probably refers to it; a CEO wouldn’t fly to China in order to resolve problems in some other growth market, so *the firm’s biggest growth market* probably refers to China.¹[todo:

¹These judgments are formalized in Grice’s Maxim of Quantity: make your contribution as informative as required, but not more so.

this is not a great example; try to find one with ambiguity that requires more than Grice to resolve.]

We can use this example to introduce some terminology:

Referring expressions include *he*, *Tim Cook*, *the firm*, *the firm's biggest growth market*. These are surface strings in the text.

Referents include TIM-COOK, APPLE, CHINA; in formal semantics, these may be viewed as objects in a model, such as a database of entities. But referents need not always be entities, as we will see.

Coreference is a property of pairs of referring expressions, which holds when they refer to the same underlying entity.

Anaphora are referring expressions whose meaning depends on another expression in context, which occurs earlier in the document or talk. **Cataphora** refer to expressions that occur later in the document, like *After she won the lottery, Susan quit her job*. **Exophora** refer to entities not defined in the linguistic context.

17.1 Forms of referring expressions

There are many possibilities for describing a referent.

- Indefinite NPs: *a visit*, *two stores*
- Definite NPs: *the capital*, *his first trip*
- Pronouns: *he*, *it*
- Demonstratives: *this chainsaw*, *that abandoned mall*
- Names: *Tim Cook*, *China*

How do you know which type of referring expression to use?

- **Language generation** requires getting this right.
You can't say: *Rob Ford apologized for "a lot of stupid things" but Rob Ford only acknowledged a video showing Rob Ford smoking what appears to be crack cocaine to demand police release it.*
- The **specific** referring expression within a type is determined by syntax and semantic constraints.

- The **type** of referring expression (pronoun, name, etc) is largely determined by the discourse.

One theory about the relationship between discourse structure and forms of referring expressions is the Givenness Hierarchy (Gundel et al., 1993). This theory is based on the **status** of the referent with respect to both the discourse and the hearer.

Type identifiable (you know what dogs are): indefinite

*I couldn't sleep, **a dog** kept me awake.*

Referential (some particular dog): indefinite *this*

*I couldn't sleep, **this dog** kept me awake.*

Uniquely identifiable definite

*I couldn't sleep, **the neighbor's dog** kept me awake.*

Familiar distal demonstrative

***That dog** next door kept me awake all night.*

Activated demonstrative

*My neighbor bought a new dog, and **that dog** kept me awake last night.*

In focus pronoun

*Her dog barks constantly. **It** kept me awake all night.*

The location of an entity in the givenness hierarchy depends (in part) on the discourse:

- *You look tired, did a dog keep you awake?*
- *We bought a dog. It keeps me up all night.*
- Referents which were recently accessed acquire *salience*, and are more likely to be near the top of the givenness hierarchy.

However, background knowledge also plays an important role.

- If a pair of speakers lives with a (single) dog, it is always at least uniquely identifiable.
- Entities may be **inferable** from the discourse:

She just bought a new bike.
The wheels are made of bamboo fiber.

Centering theory (Grosz et al., 1995) formalizes the notion of salience, by incorporating the syntactic role of each referring expression.

At each utterance U_n , we have:

- A backward-looking center $C_b(U_n)$:
the entity currently **in focus** after U_n .
- A forward-looking center $C_f(U_n)$:
an ordered list of candidates for $C_b(U_{n+1})$.
- The top choice in $C_f(U_n)$ is $C_p(U_{n+1})$

How do we order the candidates from $C_b(U_{n+1})$ to the forward-looking center?
 By syntax:

1. Subject
***Abigail** saw an elephant.*
2. Existential predicate nominal
*There is **an elephant** in the room.*
3. Direct object
*Abigail gave **a snack** to the elephant.*
4. Indirect object or oblique
*Abigail gave a snack to **the elephant**.*
5. demarcated adverbial prepositional phrase
*Inside **the zoo**, the elephant is king.*

Rule: If any element of $C_f(U_n)$ is realized by a pronoun in U_{n+1} , then $C_b(U_{n+1})$ must also be realized as a pronoun.

- Generate possible C_b and C_f for each set of reference assignments
- Filter by constraints: syntax, semantics, and centering rules
- Rank by transition orderings: continue, retain, smooth-shift, rough-shift

	$C_b(U_{n+1}) = C_b(U_n)$ or $C_b(U_n) = \emptyset$	$C_b(U_{n+1}) \neq C_b(U_n)$
$C_b(U_{n+1}) = C_p(U_{n+1})$ $C_b(U_{n+1}) \neq C_p(U_{n+1})$	Continue Retain	Smooth-shift Rough-shift

In a coherent discourse, we select transitions according to the following preferences: continue, retain, smooth-shift, rough-shift

Here's an example of how to use centering to resolve pronouns.

U_n	$C_f(U_n)$	$C_p(U_n)$	$C_b(U_n)$	transition
<i>John saw a beautiful Masi at the bike shop</i>	John, Ford, bike shop	John	\emptyset	
<i>He showed it to Bob</i>	John, Masi, Bob	John	John	Continue
<i>He showed it to Bob</i>	John, bike shop, Bob	John	John	Continue
<i>He bought it</i>	John, Masi or bike shop	John	John	Continue
<i>He bought it</i>	Bob, Masi or bike shop	Bob	Bob	Smooth-shift

- Centering theory tells us that we prefer *John* over *Bob* as the referent for *he* in U_3 , because this would be a continue transition rather than a smooth-shift.
- Centering doesn't really give us a rule for choosing *Masi* over *bike shop* in U_2 , because neither is $C_b(U_2)$. We might apply the grammatical role hierarchy since there is no other basis for this decision.

17.2 Pronouns and reference

Are all referents entities? Nope.

- *They told me that I was too ugly, but I didn't believe **it**.*
- *Alice saw Bob get angry, and I saw **it** too.*
- *They told me that I was too ugly, but **that** was a lie.*
- *Jess said she worked in security.*
*I suppose **that**'s one way to put it.*

Are all pronouns referential? Also no.

Cataphora are references to entities which are evoked *after* the reference.

When she learned what had happened, Alice took the first bus out of town.

Some pronouns have **generic** referents:

- *A good father takes care of **his** kids.*
- *I want to buy a Porsche, **they** are so fast.*
- *On the moon, **you** have to carry **your** own oxygen.*

Some pronouns don't refer to anything at all:

- Pleonastic: ***It's** raining. **It's** crazy out there.*
- Cleft: ***It's** money that she's really after.*
- Extraposition: ***It** sucks that we have to work so hard.*
- Other languages:
 - *S'il vous plaît* (literally: *if it pleases you*)
 - *Wie geht es Ihnen*

How to distinguish these from referential pronouns? Bergsma et al. (2008) propose a substitutability test.

- *You can make it in advance → You can make **them** in advance*
- *You can make it in Hollywood → You can make **them** in Hollywood*

Specifically, consider 5-gram context patterns.

*... said here Thursday that **it** is unnecessary to continue*

said	here	Thursday	that	*			
	here	Thursday	that	*	is		
		Thursday	that	*	is	unnecessary	
			that	*	is	unnecessary	to
				*	is	unnecessary	to continue

For each pattern, compute the corpus counts of five **pattern fillers**:

1. *it/its*
2. *they/them/their*
3. other pronouns *she/her/...*
4. rare words (almost always nouns)
5. all other tokens (usually nouns)

These 25 counts are converted into a feature vector, and you can train a supervised classifier.

17.3 Resolving ambiguous references

Anaphora resolution is primarily concerned with pronouns like *it, this, her*

Coreference resolution adds two additional phenomena

- **Names:** *Barack Obama, Obama, President Obama, Barry O, Nobama*
- **Nominals:** *the 44th president, the former senator from Illinois, our first African-American president*

Let's go back to our example:

Apple Inc Chief Executive Tim Cook has jetted into China for talks with government officials as **he** seeks to clear up a pile of problems in the firm's biggest growth market, from **its** contested iPad trademark to treatment of local labor. Cook is on **his** first trip to the country...

- **he** $\stackrel{?}{=}$ *Apple Inc, Tim Cook, China, talks, government officials, government, ...*
- **its** $\stackrel{?}{=}$ *the firm's biggest growth market, the firm, problems, a pile of problems, ...*
- **his** $\stackrel{?}{=}$ *Cook, local labor, its contested iPad trademark, iPad, ...*

How can we resolve these pronouns?

Semantic constraints

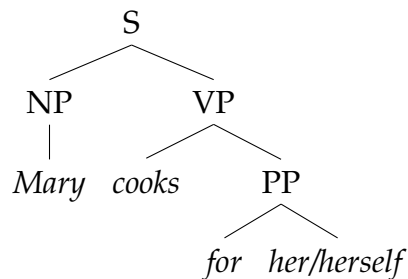
- **Number**
 - *Tim Cook has jetted in for talks with officials as **he** seeks to clear up a pile of problems...*
 - * Number(*he*) = singular
 - * Number(*officials*) = plural
 - * Number(*Tim Cook*) = singular
 - Mass noun are tricky: *New York has won the superbowl. They are the world champions.*
- **Person:** *We₁ told them₁ not to go.
- **Gender and animacy**

- *Sally met my brother. He charmed her.*
- *Sally met my brother. She charmed him.*
- *Putin brought a bottle of vodka. It was from Russia.*

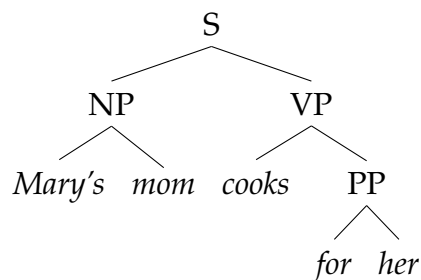
Syntactic constraints

There are general constraints on reference within sentences, which seem to generalize well across languages.

- x **c-commands** y iff the first branching node above x also dominates y .
- x **binds** y iff x and y are co-indexed and x c-commands y
- if y is not bound, it is **free**

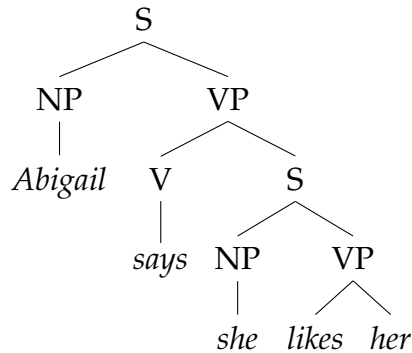


- *Mary* c-commands *her/herself*.
- *her/herself* does not c-command *Mary*.
- *her* **cannot** refer to *Mary*, because pronouns cannot refer to antecedents that c-command them.
- *herself* **must** refer to *Mary*.



- *Mary* does **not** c-commands *her*
- *Mary's mom* c-commands *her*

- *her* **can** refer to *Mary* (and we cannot use reflexive *herself* in this context, unless we are talking about Mary's mom)
- But it doesn't have to, because pronouns can be free.



Constraints have a limited domain.

- *she* can refer to *Abigail*
- *her* can also refer to *Abigail*
- But *she* and *her* cannot be coreferent.

Besides these rules, syntax also exercises preferences. See slides.

Combining the evidence

Three **types** of evidence:

- Semantic constraints
- Syntactic constraints
- Discourse/salience preferences

How do we combine them?

- **Hobbs:** Tree search + constraints

Walk back through the tree in a deterministic order, select the first referent that satisfies the constraints.

- **Centering:** ordered preferences + constraints

Apply centering theory to recover the references that give the most preferred transition sequence, subject to semantic constraints.

- **Lappin and Lease:** numerical preferences + constraints

Basically a hand-tuned linear classifier.

- -100 for each intervening sentence
- +80 for subject position
- +70 for existential emphasis, e.g. *there was a woman who...*
- +50 for accusative emphasis
- ...

- Ge, Hale, and Charniak (1999): statistical combination of four probabilities

- probability of the “Hobbs distance” between pronoun and antecedent
- probability of the pronoun given the antecedent (this considers gender and animacy)
- how well the proposed antecedent fills the pronoun’s slot in the sentence
- frequency of the proposed referent

- Raghunathan et al. (2010) describe a “multipass sieve” for coreference resolution, which applies a series of progressively relaxed matching rules.

17.4 Coreference resolution

This is a generalization of the anaphora resolution task to cover proper nouns and nominals.

- See the slides for an example.
- The coreference task comes from the information extraction community.
- Candidate spans of text for coreference are called **markables**
- In the harder versions of the coreference task, you have to identify the markables as well as their reference chains.

Coreference combines many phenomena: all the ones in anaphora resolution, plus string similarity and knowledge to get nominals.

- *unencrypted Wi-Fi networks* and *networks* have the same head word
- *Dr. King* and *Martin Luther King* can all co-refer

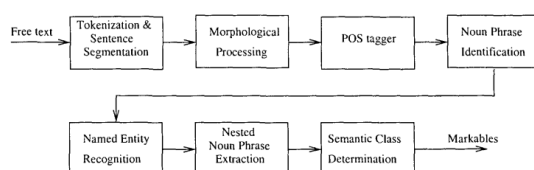
(c) Jacob Eisenstein 2014-2015. Work in progress.

- *Martin Luther King and Coretta Scott King cannot*
- **World knowledge:** e.g., *Google is a company*, companies possess *cars* but *Tuesday* doesn't.

The mention-pair model

One of the earliest end-to-end machine learning systems for coreference is from Soon et al. (2001).

- Identify markables and their features with an NLP pipeline.



- Train a classifier to predict which pairs of markables corefer. This is the **mention-pair** model.
 - For each markable, go backwards until the classifier selects an antecedent or you reach the beginning of the document.
 - No structured prediction here; each classification decision is made independently.

Learning is performed on mention pairs.

- Given the labeled chain A1-A2-A3-A4, the adjacent pairs A1-A2, A2-A3, A3-A4 are treated as positive examples.
- Negative examples are generated from NPs that occur between the adjacent pairs.
 - Suppose markables A,B,B1 appear between A1 and A2.
 - Then the negative examples are: A-A2, B-A2, B1-A2.

There are fundamental problems with mention-pair approaches.

- They fail to aggregate information across the chain.
- Must reason about transitivity to avoid incoherent chains.
- *Michelle Obama* \leftarrow *Obama* \leftarrow *Mr. Obama*

Entity-based coreference

Alternatively, we can try to learn at the entity level, using features of the entities themselves

- Number of entities detected so far
- Mention to entity ratio
- Entity to word ratio
- Number of intervening mentions between mention and linked entity
- ...

Can incorporate these by scoring entire clusterings, $\theta^\top f(x, y)$.

But how to train such a model?

One approach is an incremental perceptron. This is like a structured perceptron, but you incrementally build the structure, and you update as soon as you make a mistake.

Bell Tree, Beam Search, and Max-link Coreference The Bell Tree can represent the coreference structure. See slides.

Markov Random Field with Transitive Closure see slides

Summing over antecedent structures Durrett and Klein (2013) propose summing over reference assignments within a clustering. Let the gold standard clustering be written C^* , with C_k^* representing the cluster for document k , and $\mathcal{A}(C_k^*)$ representing the set of possible antecedents structures. Then we treat the specific antecedent structure as a latent variable, and sum over it, obtaining the regularized objective,

$$\ell(\theta) = \sum_k \log \left(\sum_{a \in \mathcal{A}(C_k^*)} p(a \mid x_k) \right) + \lambda \|\theta\| \quad (17.1)$$

$$p(a \mid x_k) \propto \exp \left(\sum_i \theta^\top f(i, a_i, x) \right). \quad (17.2)$$

Durrett and Klein (2013) augment this basic model by defining a real-valued loss function, and incorporating it into the objective. They then show that this basic framework supports a number of expressive features, which give good performance compared to prior work.

(c) Jacob Eisenstein 2014-2015. Work in progress.

17.5 Coreference evaluation

17.6 Multidocument coreference resolution

Part IV

Applications

Chapter 18

Information extraction

18.1 Entities

18.2 Relations

Knowledge-base population

Distant supervision

18.3 Events and processes

18.4 Facts, beliefs, and hypotheticals

Chapter 19

Machine translation

(These are just rough and unstructured notes right now)

Machine translation (MT) is one of the “holy grail” problems in natural language processing. Solving it would be a major advance in facilitating communication between people all over the world, and so it has received a lot of attention and funding since the early 1950s. However, it has proved incredibly challenging, and while there has been substantial progress towards usable MT systems — especially for so-called “high resource” languages like English and French — we are still far from automatically producing translations that capture the nuance and depth of human language.

19.1 The noisy channel model

Throughout the course, we’ve been working with the general formulation,

$$\hat{y} = \arg \max_{y \in \mathcal{Y}} \theta^\top f(x, y) \quad (19.1)$$

Now suppose we make \mathcal{X} equal to the set of all possible sentences in a foreign language, and \mathcal{Y} equal to the set of all possible English sentences. We can thus view translation in the same linear formalism that we’ve considered all along. Will this work?

There are two major criteria for a translation:

- **Adequacy:** The translation \hat{y} should adequately reflect the linguistic content of x . For example, if $x = \textit{Vinay le gusta Python}$, the gloss¹ $y = \textit{Vinay it like Python}$

¹A “gloss” is a word-for-word translation.

is considered adequate becomes it contains all the relevant content. The output $\mathbf{y} = \text{Vinay debugs memory leaks}$ will score poorly.

- **Fluency:** The translation $\hat{\mathbf{y}}$ should read like fluent text in the target language. By this criterion, the gloss $\mathbf{y} = \text{Vinay it like Python}$ will score poorly, and $\mathbf{y} = \text{Vinay likes Python}$ will be preferred.

	Adequate?	Fluent?
<i>Vinay it like Python</i>	yes	no
<i>Vinay debugs memory leaks</i>	no	yes
<i>Vinay likes Python</i>	yes	yes

Table 19.1: Adequacy and fluency for translations of the Spanish *Vinay le gusta Python*

An early insight in machine translation was that the scoring function for a translation can decompose across these criteria:

$$\boldsymbol{\theta}^\top \mathbf{f}(\mathbf{x}, \mathbf{y}) = \boldsymbol{\theta}_t^\top \mathbf{f}_t(\mathbf{x}, \mathbf{y}) + \boldsymbol{\theta}_\ell^\top \mathbf{f}_\ell(\mathbf{y}) \quad (19.2)$$

The features \mathbf{f}_t represent the translation model, which corresponds to the adequacy criterion; the features \mathbf{f}_ℓ represent the language model, which corresponds to the fluency criterion.

The advantage of this decomposition is that we can estimate $\boldsymbol{\theta}_\ell^\top$ from unlabeled data in the target language. Because unlabeled text data is widely available, in principle we can easily improve the fluency of our translations by estimating very high-order language models from ample unlabeled text. In this case, we can express these features as

$$\mathbf{f}_\ell(\mathbf{y}) = \bigcup_i \mathbf{1}(\mathbf{y}_{i:i+k}) \quad (19.3)$$

$$\theta_\ell(\{w_0, w_1, w_2, \dots, w_k\}) = \log p(w_k \mid w_{k-1}, w_{k-2}, \dots, w_0) \quad (19.4)$$

When estimating these probabilities, we will naturally want to apply all the smoothing tricks that we learned in Chapter 5. Note that we will also have to add padding of K “buffer” words at the beginning and end of the input.

This approach is indeed a component of the current state-of-the-art MT systems, but there is a catch: as the size of the N-gram features increases, the problem

of **decoding** — selecting the best scoring translation $\hat{\mathbf{y}}$ — becomes exponentially more difficult. We will consider this issue later. For now, just note that this formulation ensures that,

$$\boldsymbol{\theta}_\ell^\top \mathbf{f}_\ell(\mathbf{y}) = \log p(\mathbf{y}). \quad (19.5)$$

Now let's consider the translation component. If we can set

$$\boldsymbol{\theta}_t^\top \mathbf{f}_t(\mathbf{y}, \mathbf{x}) = \log p(\mathbf{x} | \mathbf{y}), \quad (19.6)$$

then the sum of these two scores yields,

$$\boldsymbol{\theta}_t^\top \mathbf{f}_t(\mathbf{y}, \mathbf{x}) + \boldsymbol{\theta}_\ell^\top \mathbf{f}_\ell(\mathbf{y}) = \log p(\mathbf{x} | \mathbf{y}) + \log p(\mathbf{y}) \quad (19.7)$$

$$= \log p(\mathbf{x}, \mathbf{y}). \quad (19.8)$$

In other words, we can obtain the translation $\hat{\mathbf{y}}$ which has the maximum joint log-likelihood $\log p(\mathbf{y}, \mathbf{x})$. We want the translation with the highest conditional probability,

$$\arg \max_{\mathbf{y}} p(\mathbf{y} | \mathbf{x}) = \arg \max_{\mathbf{y}} \frac{p(\mathbf{y}, \mathbf{x})}{p(\mathbf{x})}, \quad (19.9)$$

but since \mathbf{x} is given, we can ignore the denominator $p(\mathbf{x})$ and just select the \mathbf{y} that maximizes the joint probability.

This approach is called the **noisy channel model**, and was pioneered by researchers who were experts in cryptography. They proposed to view translation as *decoding* the output of a stochastic cipher.

- Imagine that the original text \mathbf{y} was written in English, and is modeled as drawn from a source language model $\mathbf{y} \sim P_\ell$
- The source was then stochastically encoded, according to the translation model, $\mathbf{x} | \mathbf{y} \sim P_t$.
- If we can estimate the stochastic processes P_ℓ and P_t , we can reverse the cipher and obtain the original text.

19.2 Translation modeling

Language modeling is covered in Chapter 5, so this chapter will mainly focus on the translation model, $p_t(\mathbf{x} | \mathbf{y})$. To estimate this model, we will need a parallel corpus, which contains sentences in both languages.

- Parallel corpora are often available from national and international governments. **The Hansards corpus** contains aligned English and French sentences from the Canadian parliament. **The EuroParl corpus** contains sentences for 21 languages, aligned with their English translations.
- More recent work has explored the use of web documents (Kilgariff and Grefenstette, 2003; Resnik and Smith, 2003) and crowdsourcing for MT (Zaidan and Callison-Burch, 2011).

Once a parallel corpus is obtained, we can consider how to characterize the translation model, f_t . The sets \mathcal{X} and \mathcal{Y} are far too huge for us to directly estimate the adequacy of every possible translation pair. So we need to decompose this problem into smaller units.

The **Vauquois Pyramid** is a theory of how translation should be modeled. At the lowest level, we translate individual words, but the distance here is far, because languages express ideas differently. If we can move up the triangle to syntactic structure, the distance for translation is reduced; we then need only produce target-language text from the syntactic representation, which can be as simple as reading off a tree. Further up the triangle lies semantics; translating between semantic representations should be easier still, but mapping between semantics and surface text is a difficult, unsolved problem. At the top of the triangle is **interlingua**, a semantic representation that is so generic, it is identical across all human languages. Philosophers may debate whether such a thing as interlingua is really possible (Derrida et al., 1985), but the idea of linking translation and semantic understanding is viewed by many as a grand challenge for natural language technology.

Returning to earth, the simplest decomposition of the translation model is a word-based translation: each word in the source string should be aligned to a word in the translation. In this approach, we need an alignment $\mathcal{A}(\mathbf{x}, \mathbf{y})$, which contains a list of pairs of source and target tokens. Making some independence assumptions, we can then define the translation probability as,

$$p_t(\mathbf{x}, \mathcal{A} \mid \mathbf{y}) = \prod_i p(x_i, a_i \mid y_{a_i}) \quad (19.10)$$

$$= \prod_i q(a_i \mid i, N_x, N_y) t(x_i \mid y_{a_i}) \quad (19.11)$$

Key assumptions:

- The alignment probability decomposes as $p(\cdot \mid \mathbf{x}, \mathbf{y}) = \prod_i q(a_i \mid i, N_x, N_y)$. This means that each alignment decision is independent of the others.

(c) Jacob Eisenstein 2014-2015. Work in progress.

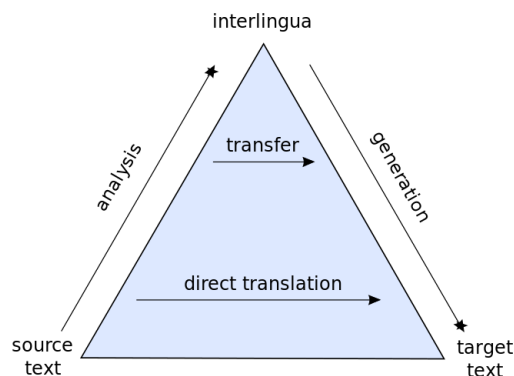


Figure 19.1: The Vauquois Pyramid (“Direct translation and transfer translation pyramid”). Licensed under Creative Commons Attribution-Share Alike 3.0 via Wikimedia Commons.)

- The translation probability decomposes as $p(\mathbf{x} \mid \mathbf{y}, \mathbf{y}) = \prod_i t(x_i | y_{a_i})$. We are doing word-based translation only, ignoring context.

A series of translation models with increasingly weak independence assumptions was produced by researchers at IBM in the 1980s and 1990s, known as IBM Models 1-6 (Och and Ney, 2003). IBM model 1 makes the strongest independence assumption:

$$q(a_i \mid i, N_x, N_y) = \frac{1}{N_y} \quad (19.12)$$

Let’s consider how to translate with IBM model 1. The key idea is to treat the alignment as a **hidden variable**. If we knew the alignment, we could easily estimate a translation model, and we could find the optimal translation as

$$\hat{\mathbf{y}} = \arg \max_{\mathbf{y}} p(\mathbf{x}, \mathbf{y}) \quad (19.13)$$

$$= \arg \max_{\mathbf{y}} \sum_{\mathcal{A}} p(\mathbf{x}, \mathbf{y}, \mathcal{A}) \quad (19.14)$$

$$= \arg \max_{\mathbf{y}} p_{\ell}(\mathbf{y}) \sum_{\mathcal{A}} p_t(\mathbf{x}, \mathcal{A} \mid \mathbf{y}) \quad (19.15)$$

$$\approx \arg \max_{\mathbf{y}} p_{\ell}(\mathbf{y}) \max_{\mathcal{A}} p_t(\mathbf{x}, \mathcal{A} \mid \mathbf{y}) \quad (19.16)$$

Conversely, if we knew the translation model, we could estimate beliefs about each alignment decision $p(a_i \mid \mathbf{x}, \mathbf{y}) \propto q(a_i) p_t(x_i \mid y_{a_i})$.

We therefore have a classic chicken-and-egg problem, which we can solve using the iterative expectation-maximization (EM) algorithm.

E-step Update beliefs about word alignment,

$$p(a_i \mid \mathbf{x}, \mathbf{y}) \propto q(a_i) p_t(\mathbf{x}_i \mid \mathbf{y}_{a_i}) \quad (19.17)$$

M-step Update translation model,

$$\theta_{u \rightarrow v} = \log \sum_i \sum_j p(a_i = j) \delta(\mathbf{y}_j = u) \delta(\mathbf{x}_i = v) - \log \sum_i \sum_j p(a_j = j) \delta(\mathbf{y}_j = u) \quad (19.18)$$

19.3 Example for IBM Model 1

- Translation probabilities

	<i>the</i>	<i>house</i>
<i>la</i>	0.4	0.1
<i>maison</i>	0.1	0.6

- Alignments

	$P(\mathbf{f}, \mathbf{a} \mid \mathbf{e})$	$P(\mathbf{a} \mid \mathbf{f}, \mathbf{e})$
<i>the</i> \rightarrow <i>la</i> , <i>house</i> \rightarrow <i>la</i>	$\lambda \times 0.4 \times 0.1 = 0.04\lambda$	0.11
<i>the</i> \rightarrow <i>la</i> , <i>house</i> \rightarrow <i>maison</i>	$\lambda \times 0.4 \times 0.6 = 0.24\lambda$	0.69
<i>the</i> \rightarrow <i>maison</i> , <i>house</i> \rightarrow <i>la</i>	$\lambda \times 0.1 \times 0.6 = 0.06\lambda$	0.17
<i>the</i> \rightarrow <i>maison</i> , <i>house</i> \rightarrow <i>maison</i>	$\lambda \times 0.1 \times 0.1 = 0.01\lambda$	0.03

- Counts

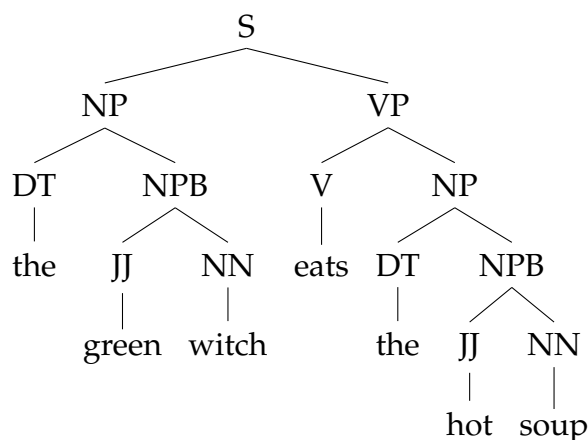
	<i>the</i>	<i>house</i>
<i>la</i>	$0.11 + 0.69 = 0.8$	$0.11 + 0.17 = 0.28$
<i>maison</i>	$0.17 + 0.03 = 0.2$	$0.69 + 0.03 = 0.72$

Then we add up the counts across all the examples and update the translation probabilities

19.4 Syntactic MT

Consider the English sentence, *The green witch eats the hot soup*.

(c) Jacob Eisenstein 2014-2015. Work in progress.



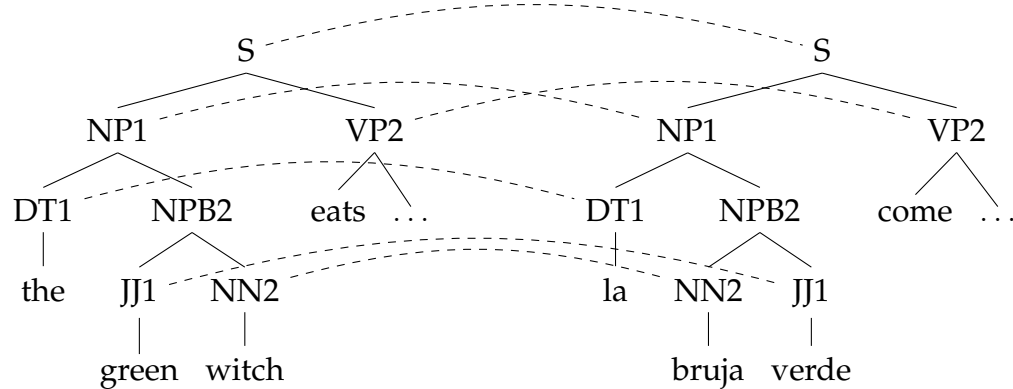
Where NPB is a “bare NP,” without the determiner. We might get this non-terminal from binarizing a CFG.

We can view the CFG as a process for **generating** English sentences.

Synchronous CFGs are a generalization of CFGs. They generate text in two different languages simultaneously. Each RHS has two components, one for each language. Subscripts show the mapping between non-terminals in the RHS. For example:

$$\begin{array}{ll}
 S \rightarrow NP_1 VP_2, & NP_1 VP_2 \\
 VP \rightarrow V_1 NP_2, & V_1 NP_2 \\
 NP \rightarrow DT_1 NPB_2, & DT_1 NPB_2 \\
 NPB \rightarrow JJ_1 NPB_2, & NPB_2 JJ_1
 \end{array}$$

The key production is the fourth one, which handles the re-ordering of adjectives and nouns. Let’s use this SCFG to generate the English and Spanish versions of this sentence.



- On the slides there is another example, in Japanese. Since Japanese is a SOV language (subject-object-verb), we need a production: $VP \rightarrow V_1 NP_2, NP_2 V_1$.
- As with CFGs, we can attach a probability to each production, and compute the joint probability of the derivation and the text as the product of these productions.

Binarization

Let's define a rank- n CFG as a grammar with at most n elements on a right-hand side.

- CFGs can always be binarized.
 - e.g. $NP \rightarrow DT [JJ NN]$ becomes

$$NP \rightarrow DT NPB$$

$$NPB \rightarrow JJ NN$$
 - Therefore, the set of languages that can be defined by a 2-CFG is identical to the set that can be defined by 3-CFG, 4-CFG, etc...

- What about SCFGs?

- Rank 3:

$$\begin{array}{ll}
 A \rightarrow B [C D], & [C D] B \\
 A \rightarrow B V, & V B \\
 V \rightarrow C D, & C D
 \end{array}$$

Yes, we can. 2-SCFG = 3-SCFG.

(c) Jacob Eisenstein 2014-2015. Work in progress.

– Rank 4:

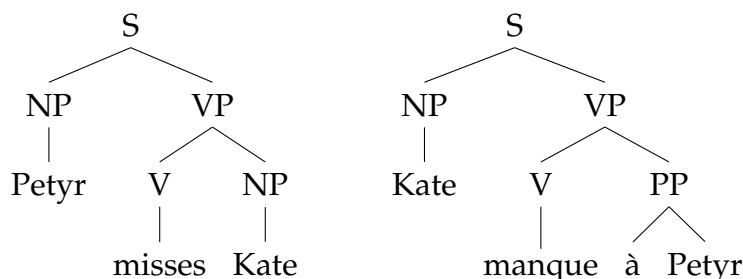
$$\begin{array}{ll}
 A \rightarrow B C D E, & C E B D \\
 A \rightarrow [B C] D E, & [C E B] D \\
 A \rightarrow B [C D] E, & [C E B D] \\
 A \rightarrow B C [D E], & C [E B D]
 \end{array}$$

In each chunk that we might want to replace in the first language, we have one or more intervening symbols in the second language. Therefore, 3-SCFG \subsetneq 4-SCFG.

- The subset of 2-SCFG = 3-SCFG is equivalently called **inversion transduction grammar**. The notation is slightly different, we write $A \rightarrow [B C]$ when the order is preserved and $A \rightarrow \langle B C \rangle$ when it is inverted.

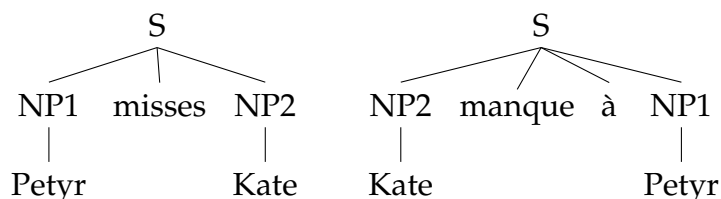
No raising or lowering

SCFGs can only reorder sibling nodes. Is that enough? Not always.



SCFGs cannot swap the subject and object, because they aren't siblings in the original grammar.

We could solve this by changing the grammar,



By including the verb *misses/manque à* directly into the rule, we ensure that it doesn't apply to other verbs.

With other syntactic translation models (synchronous tree substitution grammar or tree adjoining grammars), this case can be handled without flattening.

19.5 Algorithms for SCFGs

Translation

In principle, translation in SCFGs is nearly identical to parsing. Suppose we have the Spanish phrase *la razón principal*, and the synchronous grammar

$NP \rightarrow D\ NPB,$	$D\ NPB$	1.0
$NPB \rightarrow N_1\ J_2,$	$J_2\ N_1$	0.8
$NPB \rightarrow N_1\ N_2,$	$N_1\ N_2$	0.2
$D \rightarrow la,$	the	0.5
$N \rightarrow razon,$	$reason$	0.5
$N \rightarrow principal,$	$principal$	0.5
$J \rightarrow principal,$	$main$	1.0

Now we can apply CKY, building the translation on the English side. We should get two possible translations, *the reason principal* ($p(e, f, \tau) = 0.05$) and *the main reason* ($p(e, f, \tau) = 0.4$).

What is the complexity of translation with binarizable SCFGs? It's just like CFG parsing: $\mathcal{O}(n^3)$.

Bitext parsing

To learn a translation model, we might need to synchronously parse the **bitext**: both the source and target side language.

We can do this with a dynamic program.

Assuming we are dealing with 2-SCFG or 3-SCFG, here's what we need to keep track of:

- The non-terminals that we have derived
- Their spans in the source language (start and end)
- Their spans in the target language (start and end)

Suppose we are given spans $\langle i, j \rangle$ in the source and $\langle i', j' \rangle$ in the target. Then we are looking for split points k and k' and a production that can derive the subspans $\langle i, k \rangle$, $\langle k, j \rangle$ and $\langle i', k' \rangle$, $\langle k', j' \rangle$.

What is the space complexity of bitext parsing? $\mathcal{O}(|S|n^4)$, where $|S|$ is the number of non-terminals.

What is the time complexity of bitext parsing? $\mathcal{O}(|R|n^6)$, where $|R|$ is the number of production rules.

Specifically, we have the recurrence

$$\begin{aligned} \psi(X, i, j, i', j') = & \max_{k, k', A, B} P(S \rightarrow A B, A B) \otimes \psi(A, i, k, i', k') \otimes \psi(B, k, j, k', j') \\ & \oplus P(S \rightarrow A B, B A) \otimes \psi(A, i, k, k', j') \otimes \psi(B, k, j, i', k') \end{aligned}$$

Note: in general, bitext parsing is exponential in the rank of the SCFG (unless $P = NP$).

Intersection with language model

For fluent translations, we typically want to multiply in the language model probability on the target side.

- This (usually) corresponds **intersection** of an SCFG with a finite state machine.
- Sidenote: **what about context-free language models?**
 - $A = \{a^m b^m c^n\}$
 - $B = \{a^m b^n c^n\}$
 - $A \cap B = \{a^n b^n c^n\}$, not a CFL!
 - CFLs are not closed under intersection.
 - Determining if $s \in A \cap B$ is in PSPACE
- There are exact dynamic programming algorithms for intersecting an SCFG and an FSA, but they are very slow. One solution is **cube pruning**.
- We can equivalently view this as an ILP

$$\begin{aligned} \min. \quad & \sum_v \theta_v y_v + \sum_e \theta_e y_e + \sum_{\langle v, w \rangle \in \mathcal{B}} \theta(v, w) y(v, w) \\ \text{s.t.} \quad & C0 : y_v, y_e \text{ form a derivation} \\ & C1 : y_v = \sum_{w: \langle w, v \rangle \in \mathcal{B}} y(w, v) \\ & C2 : y_v = \sum_{w: \langle v, w \rangle \in \mathcal{B}} y(v, w) \end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Here y_e and y_v are indicator variables that define what words and hyper-edges appear in the derivation.
- We can solve this optimization with Lagrangian relaxation.
 - Replace the outgoing constraints $C2$ with multipliers $u(v)$
 - At first, $u(v) = 0, \forall v$
 - Without the outgoing constraints, we can optimize efficiently
 - If the outgoing constraints happen to be met, we are done
 - Otherwise, update $u(v)$ and try again.
- Lagrangian relaxation finds the exact solution 97% of the time, is many times faster than ILP.

Part V

Learning

Chapter 20

Semi-supervised learning

So far we have focused on learning a classifier — typically represented by a set of weights θ — from a set of labeled examples $\{(x_i, y_i)\}_{i=1}^{\ell}$. As we've seen, it's possible to formulate structured prediction tasks such as parsing in this same framework. But what if you don't have those labeled examples for the domain or task that you want to solve?

- You can use some other labeled data and hope it works.
This rarely works well.
- You can label data yourself.
This is a lot of work.

This kind of thing happens all the time (especially in class projects). And labeled data is really expensive:

- **The Switchboard corpus** contains phoneme annotations of telephone conversations, e.g.

film → F IH_N UH_GL_N M
be all → BCL B IY IY_TR AO_TR AO L_DL

This took 400 hours of annotation time per hour of speech.

- **The Penn Chinese Treebank** is a set of CFG annotations for Chinese. It took 2 years to get 4000 sentences annotated.

20.1 Learning with less annotation effort

We can think of our annotated data as a *sample* from some underlying distribution.

$$\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{\ell} \sim \mathcal{D} \quad (20.1)$$

This allows us to formulate various learning scenarios:

Semisupervised learning

- $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{\ell} \sim \mathcal{D}$: labeled examples
- $\{(\mathbf{x}_i)\}_{i=\ell+1}^{\ell+u}$: unlabeled examples
- often $u \gg \ell$

We’ve already seen an example of semi-supervised learning in document classification, when we applied expectation maximization to impute labels of unlabeled documents. Today we will see some approaches that tend to work better than EM.

Active learning is closely related to semi-supervised learning, but you can now query the labels for a few examples while learning. Your goal is to select these examples carefully, so that you get the best accuracy from a fixed number of examples, or so that you get to a certain level of accuracy with as few examples as possible.

Domain adaptation

- $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{\ell_S} \sim \mathcal{D}_S$: labeled examples in *source* domain
- Some *target domain* \mathcal{D}_T . You may have a small amount of labeled data in the target domain (“supervised domain adaptation”) or not (“unsupervised domain adaptation”).

The prototypical example for domain adaptation is product reviews: you have annotated reviews of coffee machines, but you want to train a classifier for reviews of bicycles. Another example is that you have a POS tagger for news genre text, but you want one for social media.

- Suppose that $p_T(X) = p_S(X)$, i.e. the distribution over observed data is the same in both domains, but $p_T(Y|X) \neq p_S(Y|X)$, i.e. the conditional distribution over labels is different. This setting is sometimes called **transfer learning** or **multitask learning**.

(c) Jacob Eisenstein 2014-2015. Work in progress.

For example, you have labeled data for part-of-speech tagging, but you want to train a system for named entity recognition. Or, you have a classifier for detecting mentions of people and you want one for detecting mentions of places.

20.2 Why would unlabeled data help?

Suppose you want to do sentiment analysis in French. I give you two labeled examples:

- ☺ émouvant avec grâce et **style**
- ☹ fastidieusement inauthentique et **banale**

You have a bunch of unlabeled examples too:

1. pleine de **style** et d'**intrigue**
2. la **banalité** n'est dépassée que par sa **prétention**
3. **prétentieux**, de la première minute au rideau final
4. imprégné d'un air d'**intrigue**

What can we do? If we just learn from the labeled data, we might decide that *style* is positive and that **banale** is negative. This isn't much. However, we can propagate this information to the unlabeled data, and potentially learn more.

- If we are confident about *style* being positive, then we can guess that example 1 is also positive.
- That suggests that *intrigue* is also positive.
- We can then propagate this information to example 4, and learn more.
- Similarly, we can propagate from the labeled data to example 2, which we guess to be negative. This suggests that *pretention* is also negative, which we propagate to example 3.

What happened here?

- Instances 1 and 2 were “similar” to our labeled examples for positivity and negativity, respectively. We used them to expand those concepts, which allowed us to correctly label instances 3 and 4, which didn’t share any important features with our original labeled data.
- We made a key assumption: that similar instances will have similar labels. Is this a strong assumption? Keep this question in mind.
- In this case, we defined similarity in terms of sharing some key words (non-stopwords).
- To see how this can help conceptually, think about similarity just in terms of 1D space. If you have only the two labeled instances, your decision boundary should be right in between. (Do you remember what criterion justifies this choice?) But if you have a bunch of unlabeled instances, you might want to draw this boundary in a different place.
- Let’s see how we can operationalize this idea in an algorithm.

Semi-supervised learning with EM

We’ve already seen one way to do this: use expectation-maximization to marginalize over the labels of the unseen data. So we are maximizing

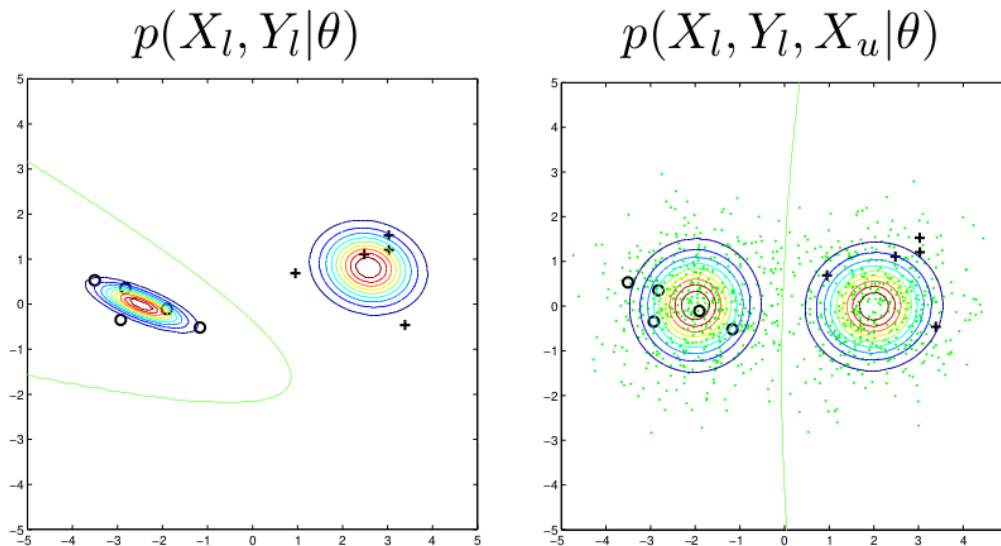
$$p(X^\ell, Y^\ell, X^U) = p(X^\ell, Y^\ell) \sum_{Y^U} p(X^U, Y^U) \quad (20.2)$$

We did this by

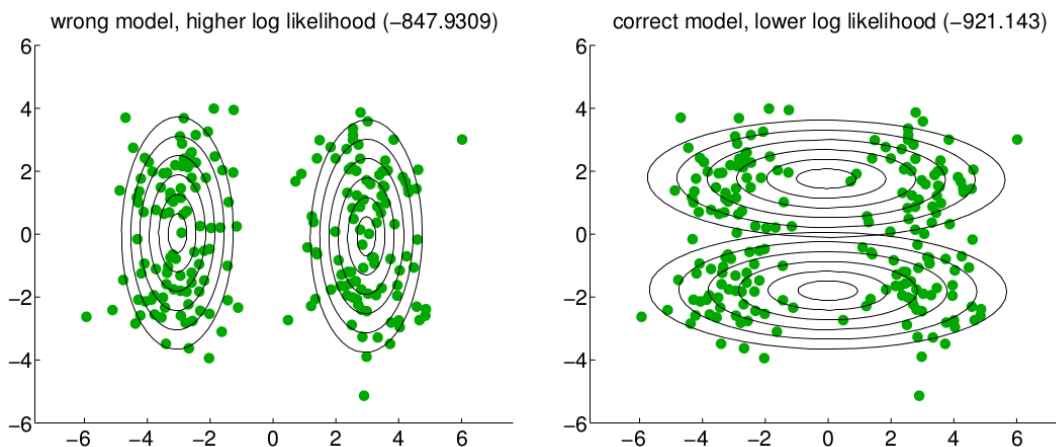
- **E.** fitting a distribution $Q(y_i)$ for all unlabeled i ,
- **M.** maximizing the expected likelihood under this distribution.

You can see why this can work in an example:

(c) Jacob Eisenstein 2014-2015. Work in progress.



We get a much more reasonable decision boundary. However, things can also go wrong:



The correct model has a lower log-likelihood than the incorrect model. The basic problem is that our model is wrong. The label is related to the observations, but not in the simplistic, Gaussian way that we had assumed. We discussed a heuristic to try to deal with this problem: downweighting the contribution of the unseen data to the likelihood function.

(c) Jacob Eisenstein 2014-2015. Work in progress.

Bootstrapping

EM is sort of like self-training or **bootstrapping**: we use our current model to estimate $Q(y_i)$, and then update the model as if $Q(y_i)$ is correct.

- The probabilistic nature of this is nice, but it limits us to relatively weak classifiers.
- If we are willing to give up on probability, we can bootstrap **any** classifier.

We can try this first using one nearest-neighbor (see slides). Like EM, it can work, but it can also fail. (The failure modes are different though; can you characterize the difference?)

Co-training

“Folk wisdom” about when bootstrapping works:

- Better for generative models (e.g., Naive Bayes) than for discriminative models (e.g., perceptron)
- Better when the Naive Bayes assumption is stronger.
 - Suppose we want to classify NEs as PERSON or LOCATION
 - Features: string and context
 - * *located on Peachtree Street*
 - * *Dr. Walker said ...*

$$\begin{aligned} P(X_1 = \text{street}, X_2 = \text{on} \mid Y_1 = \text{LOC}) \\ \approx P(X_1 = \text{street} \mid Y_1 = \text{LOC})P(X_2 = \text{on} \mid Y_1 = \text{LOC}) \end{aligned}$$

Cotraining makes the bootstrapping assumptions explicit.

- Assume two, **conditionally independent**, views of a problem.
- Assume each view is sufficient to do good classification.

Sketch of learning algorithm:

- On labeled data, minimize error.
- On unlabeled data, **constrain** the models from different views to agree with each other.

Co-training example See the slides for an animated version of this. Assume we want to do named entity classification: determine whether an NE is a Location or Person. We have two views: the name itself, and its context.

	$\mathbf{x}^{(1)}$	$\mathbf{x}^{(2)}$	y
1.	Peachtree Street	located on	LOC
2.	Dr. Walker	said	PER
3.	Zanzibar	located in	? \rightarrow LOC
4.	Zanzibar	flew to	? \rightarrow LOC
5.	Dr. Robert	recommended	? \rightarrow PER
6.	Oprah	recommended	? \rightarrow PER

Algorithm

- Use classifier 1 to label example 5.
- Use classifier 2 to label example 3.
- Retrain both classifiers, using newly labeled data.
- Use classifier 1 to label example 4.
- Use classifier 2 to label example 6.

Multiview Learning is another approach in this style.

- Co-training treats the output of each view's classifier as a labeled instance for the other view.
- In multiview learning, we add a **co-regularizer** that penalizes disagreement between the views on the unlabeled instances.
- This allows us to define a single objective function. In the case of two-view linear regression, the function is

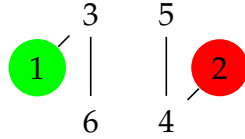
$$\min_{\mathbf{w}, \mathbf{v}} \sum_i^L (y_i - \mathbf{w}^\top \mathbf{x}_i^{(1)})^2 + (y_i - \mathbf{v}^\top \mathbf{x}_i^{(2)})^2 + \lambda_1 \|\mathbf{w}\|^2 + \lambda_1 \|\mathbf{v}\|^2 + \lambda_2 \sum_{i=L+1}^{L+U} (\mathbf{w}^\top \mathbf{x}_i^{(1)} - \mathbf{v}^\top \mathbf{x}_i^{(2)})^2$$

The only difference from standard regression is the co-regularizer, which penalizes disagreement on the unlabeled data.

- An early version of this idea is **co-boosting** (Collins and Singer, 1999), where each view is a boosting classifier, and features are added incrementally to each view.

Graph-based approaches

Let's go back to sentiment analysis in French.



We can view this data as a **graph**, with edges between similar instances. Unlabeled instances propagate information through the graph.

Where does the graph come from?

- Sometimes there is a natural similarity metric (time, position in the document).
- Otherwise, we can compute similarity from features. If the features are Gaussian, we could say:

$$\text{sim}(i, j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

If the features are discrete, we might use KL-divergence.

- Then we add an edge between i and j when $\text{sim}(i, j) > \tau$

Given a graph with edge weights s_{ij} , we can formulate semi-supervised learning as an optimization problem:

$$\begin{aligned} y_i &\in \{0, 1\} \\ \text{Fix } Y_l &= \{y_1, y_2, \dots, y_\ell\} \\ \text{Solve for } Y_u &= \{y_{\ell+1}, \dots, y_{\ell+m}\} \\ \min_{Y_u} \sum_{i,j} s_{ij} (y_i - y_j)^2 \end{aligned}$$

- This looks like a combinatorial problem. Specifically, it looks like (binary) integer linear programming, which is NP-complete.
- But assuming $s_{ij} \geq 0$, this specific problem can be reformulated as maximum-flow, with polynomial time solutions.

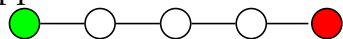
(c) Jacob Eisenstein 2014-2015. Work in progress.

Rao and Ravichandran (2009) apply this idea to expanding polarity lexicons.

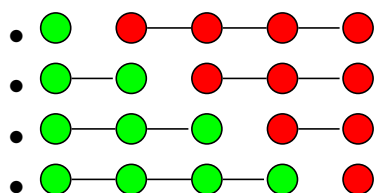
- Nodes are words
- Edges are wordnet relations
- They label a few nodes for sentiment polarity, and propagate those labels to other parts of the graph.
- However, they use a slightly modified version of the mincut idea: randomized min-cut (Blum et al., 2004).

Randomized min-cut

Suppose we have this initial graph:



What is the solution? Actually, the following solutions are all equivalent:



Another problem with mincuts is that it doesn't distinguish high-confidence and low-confidence predictions. Both of these problems can be dealt with by randomization:

- Add random noise to adjacency matrix.
- Rerun mincuts multiple times.
- Deduce the final classification by voting.

Label propagation

A related approach is **label propagation** (Zhu and Ghahramani, 2002), which Rao and Ravichandran also consider. The basic idea is that we relax y_i from an integer $\{0, 1\}$ to a real number \mathbb{R} . Then we solve the optimization problem,

$$\min_Y \sum_{i,j} s_{ij} (y_i - y_j)^2$$

s.t. Y_L is clamped to initial values

The advantages are:

- a unique global optimum
- a natural notion of confidence: distance of y_i from 0.5

Let's look at the objective:

$$\begin{aligned}
 J &= \frac{1}{2} \sum_{i,j} s_{ij} (y_i - y_j)^2 \\
 &= \frac{1}{2} \sum_{i,j} s_{ij} (y_i^2 + y_j^2 - 2y_i y_j) \\
 &= \sum_i y_i^2 \sum_j s_{i,j} - \sum_{i,j} s_{ij} y_i y_j \\
 &= \mathbf{y}^\top \mathbf{D} \mathbf{y} - \mathbf{y}^\top \mathbf{S} \mathbf{y} \\
 &= \mathbf{y}^\top \mathbf{L} \mathbf{y}
 \end{aligned}$$

We have introduced three matrices

- Let \mathbf{S} be the $n \times n$ similarity matrix.
- Let \mathbf{D} be the **degree matrix**, $d_{ii} = \sum_j s_{ij}$. \mathbf{D} is diagonal.
- Let \mathbf{L} be the unnormalized **graph Laplacian** $\mathbf{L} = \mathbf{D} - \mathbf{S}$
- So we want to minimize $\mathbf{y}^\top \mathbf{L} \mathbf{y}$ with respect to \mathbf{y}_u , the labels of the unannotated instances.

In principle, this is easily solveable:

- Partition the Laplacian $\mathbf{L} = \begin{bmatrix} \mathbf{L}_{\ell\ell} & \mathbf{L}_{\ell u} \\ \mathbf{L}_{u\ell} & \mathbf{L}_{uu} \end{bmatrix}$
- Then the closed form solution is $\mathbf{y}_u = -\mathbf{L}_{uu}^{-1} \mathbf{L}_{u\ell} \mathbf{y}_\ell$
- This is great ... if we can invert \mathbf{L}_{uu} .

In practice, $\mathbf{L}_{u,u}$ is huge, so we can't invert it unless it has special structure. Zhu and Ghahramani (2002) propose an iterative solution called **label propagation**.

- Let $\mathbf{T}_{ij} = \frac{s_{ij}}{\sum_k s_{kj}}$, row-normalizing \mathbf{S} .
- Let \mathbf{Y} be an $n \times C$ matrix of labels, where C is the number of classes.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Until tired,
 - Set $\mathbf{Y} = \mathbf{T}\mathbf{Y}$
 - Row-normalize \mathbf{Y}
 - Clamp the seed examples in \mathbf{Y} to their original values
- There's a flavor of EM here, with \mathbf{Y} representing our belief $q_i(y_i)$. But there's no M-step in which we update model parameters. That's because we're in a graph-based framework, and we're assuming the graph is correct.

Both mincut and label propagation are **transductive** learning algorithms: they learn jointly over the training and test data. This is fine in some settings, but not if you want to train a system and then apply it to new test data later — you'd have to retrain it all over again.

Manifold regularization (Belkin et al., 2006) addresses this issue, by learning functions that are smooth on the “graph manifold.” In practice, this means that they give similar labels to nearby datapoints in the unlabeled data. Suppose we are interested in learning a classification function f . Then we can optimize:

$$\arg \min_f \frac{1}{\ell} \sum_i \ell(f(\mathbf{x}_i), y_i) + \lambda_1 \|f\|^2 + \lambda_2 \sum_{i,j} s_{ij} (f(\mathbf{x}_i) - f(\mathbf{x}_j))^2$$

- The first term corresponds to the loss on the labeled training data; we can use any convex loss functions, such as logistic or hinge loss.
- The second term corresponds to the smoothness, akin to regularizing the weights in a linear classifier.
- The third term penalizes making different predictions for similar instances in the unlabeled data

The representer theorem guarantees that we can solve for f as long as ℓ is convex. We can then apply f to any new unlabeled test data.

20.3 Domain adaptation

In domain adaptation, we have a lot of labeled data, but it's in the wrong domain. Some features will be shared across domains. For example, if we are classifying movies or toasters, *good* is a good word, and *sucks* is a bad word. But as we've seen, real review text is usually more subtle. What about a word like *unpredictable*? This is a good word for a movie, but not such a good word for a kitchen appliance.

Supervised domain adaptation

In supervised domain adaptation (transfer learning), we have:

- Lots of labeled data in a “source” domain, $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{\ell_S} \sim \mathcal{D}_S$ (e.g., reviews of restaurants)
- A little labeled data in a “target” domain, $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^{\ell_T} \sim \mathcal{D}_T$ (e.g., reviews of chess stores)

Here are some (surprisingly-competitive) baselines (see slides)

- Source-only: train on the source data, apply it to the target data.
- Target-only: forget the source data, just train on the limited target data.
- Big blob: merge the source and target data into a single training set. Optionally downweight the source data.
- Prediction: train a classifier on the source data, use its prediction as a feature in the target data.
- Interpolation: train two classifiers, combine their outputs

Here are two less-obvious approaches:

Priors :

Train a (logistic-regression) classifier on the source data. Treat its weights as the priors on the target data, and regularize towards these weights rather than towards zero (Chelba and Acero 2004).

Feature augmentation Create **copies** of each feature, for each domain and for the cross-domain setting.

- The copies fire in the appropriate domains, and the learning algorithm decides whether a feature is general or domain-specific.
- For example, suppose we have domains for Appliances and Movies, and features *outstanding* and *sturdy*. We replicate the features, obtaining

$$\begin{aligned} &\langle \textit{outstanding}, \text{APP.} \rangle, \langle \textit{outstanding}, \text{MOV.} \rangle, \langle \textit{outstanding}, \text{ALL} \rangle \\ &\langle \textit{sturdy}, \text{APP.} \rangle, \langle \textit{sturdy}, \text{MOV.} \rangle, \langle \textit{sturdy}, \text{ALL} \rangle \end{aligned}$$

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Ideally, we will learn a positive weight for $\langle \textit{outstanding}, \text{ALL} \rangle$, because the feature works in both domains, and a small weight for the domain-specific copies of the *outstanding* feature.
- We will also learn a positive weight for $\langle \textit{sturdy}, \text{APP} \rangle$, because the feature works only in the Appliance domain.

See slides for a diagram of how this works.

Unsupervised domain adaptation

Without labeled data in the target domain, can we learn anything? If the source and target domain are somewhat related, then we can. A very popular approach is structural correspondence learning (SCL) (Blitzer et al., 2007).

- Suppose there are a few words that are good predictors in both domains; we'll call these **pivot features**
- Pivot features can be selected by finding words that are
 - Popular in both domains
 - High mutual-information with the label in the source domain
- The label is unknown in the target domain, so we can't learn to predict it. Instead we'll predict the pivots. We train a linear classifier for each pivot, obtaining weights θ_n for pivot n .
- For example, we can learn that the domain-specific feature *fast-multicore* is a good predictor of the pivot *excellent*.
- We can horizontally concatenate the pivot predictor weights, forming

$$\Theta = [\theta_1, \theta_2, \dots, \theta_N] \quad (20.3)$$

- The matrix Θ is large, and contains redundant information (since many pivots are closely related to each other). We factor $\Theta \approx USV^T$ using singular value decomposition (SVD).
- We use U to **project** features from both domains into a shared space, $U^\top x$.
- We then learn to predict the label in the source domain, using the augmented instance $\langle x, U^\top x \rangle$. In U contains meaningful correspondences between the domains, then the weights learned on these features will work for the target domain instances too.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- This idea yields substantial improvements in adapting sentiment classifiers across product domains, e.g., books, movies, and appliances (Blitzer et al., 2007).

See the slides for a graphical explanation of these ideas, with slightly different notation.

20.4 Other learning settings

There are many other settings in which we learn from something other than in-domain labeled data:

- **Active learning.** The model can query the annotator for labels (see above)
- **Feature labeling.** Annotators label *features* rather than instances. For example, you provide a list of five prototype words for each POS tag (Haghighi and Klein, 2006).
- **Feature expectations.** Learn from *constraints* on feature-label relationships; for example, the word “the” is a determiner at least 90% of the time. In EMNLP 2013, this idea was applied to multilingual learning (which I’ll discuss in the final lecture). The basic idea of this paper is to align words between sentences and insist that aligned words have the same tag most of the time.
- **Multi-instance learning.** The learner gets a “bag” of instances, and a label. If the label is positive, then at least one instance in the bag is positive, but you don’t know which one.

This idea is often related to **distant supervision**. The learner gets a label indicating that there is a relationship, such as BORN-IN(OBAMA, HAWAII), and a set of instances containing sentences that mention the two arguments, *Obama* and *Hawaii*. Many of these sentences do not actually instantiate the desired relation (e.g., *Obama visited Hawaii in 2008...*), but we assume that at least one does, and we must learn from this.

Chapter 21

Beyond linear models

21.1 Representation learning

21.2 Convolutional neural networks

21.3 Recursive neural networks

21.4 Encoder-decoder models

Bibliography

- Abney, S. P. and Johnson, M. (1991). Memory requirements and local ambiguities of parsing strategies. *Journal of Psycholinguistic Research*, 20(3):233–250.
- Akaike, H. (1974). A new look at the statistical model identification. *Automatic Control, IEEE Transactions on*, 19(6):716–723.
- Allauzen, C., Riley, M., and Schalkwyk, J. (2009). A generalized composition algorithm for weighted finite-state transducers. In *Proceedings of the International Speech Communication Association (INTERSPEECH)*.
- Anandkumar, A., Ge, R., Hsu, D., Kakade, S. M., and Telgarsky, M. (2014). Tensor decompositions for learning latent variable models. *The Journal of Machine Learning Research*, 15(1):2773–2832.
- Arora, S., Ge, R., Halpern, Y., Mimno, D., Moitra, A., Sontag, D., Wu, Y., and Zhu, M. (2013). A practical algorithm for topic modeling with provable guarantees. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 280–288.
- Belkin, M., Niyogi, P., and Sindhvani, V. (2006). Manifold regularization: A geometric framework for learning from labeled and unlabeled examples. *The Journal of Machine Learning Research*, 7:2399–2434.
- Bender, E. M. (2013). *Linguistic Fundamentals for Natural Language Processing: 100 Essentials from Morphology and Syntax*, volume 6 of *Synthesis Lectures on Human Language Technologies*. Morgan & Claypool Publishers.
- Bengio, Y., Ducharme, R., Vincent, P., and Janvin, C. (2003). A neural probabilistic language model. *The Journal of Machine Learning Research*, 3:1137–1155.

- Berg-Kirkpatrick, T., Bouchard-Côté, A., DeNero, J., and Klein, D. (2010). Painless unsupervised learning with features. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 582–590, Los Angeles, CA.
- Berger, A. L., Pietra, V. J. D., and Pietra, S. A. D. (1996). A maximum entropy approach to natural language processing. *Computational linguistics*, 22(1):39–71.
- Bergsma, S., Lin, D., and Goebel, R. (2008). Distributional identification of non-referential pronouns. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 10–18, Columbus, OH.
- Bikel, D. M. (2004). Intricacies of Collins’ parsing model. *Computational Linguistics*, 30(4):479–511.
- Bird, S., Klein, E., and Loper, E. (2009). *Natural language processing with Python*. O’Reilly Media, California.
- Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4):77–84.
- Blei, D. M. (2014). Build, compute, critique, repeat: Data analysis with latent variable models. *Annual Review of Statistics and Its Application*, 1:203–232.
- Blitzer, J., Dredze, M., and Pereira, F. (2007). Biographies, bollywood, boom-boxes and blenders: Domain adaptation for sentiment classification. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 440–447, Prague.
- Blum, A., Lafferty, J., Rwebangira, M. R., and Reddy, R. (2004). Semi-supervised learning using randomized mincuts. In *Proceedings of the International Conference on Machine Learning (ICML)*, page 13.
- Bottou, L. (1998). Online learning and stochastic approximations. *On-line learning in neural networks*, 17:9.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, New York, NY, USA.
- Brants, T. (2000). Tnt: a statistical part-of-speech tagger. In *Proceedings of the sixth conference on Applied natural language processing*, pages 224–231. Association for Computational Linguistics.
- Briscoe, T. (2011). Introduction to formal semantics for natural language.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- (c) Jacob Eisenstein 2014-2015. Work in progress.

- Church, K. W. (2000). Empirical estimates of adaptation: the chance of two Norie-gas is closer to $p/2$ than p^2 . In *Proceedings of the 18th conference on Computational linguistics-Volume 1*, pages 180–186.
- Collins, M. (1997). Three generative, lexicalised models for statistical parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 16–23.
- Collins, M. (2002). Discriminative training methods for hidden markov models: theory and experiments with perceptron algorithms. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1–8.
- Collins, M. (2003). Head-driven statistical models for natural language parsing. *Computational linguistics*, 29(4):589–637.
- Collins, M. (2013). Notes on natural language processing. <http://www.cs.columbia.edu/~mcollins/notes-spring2013.html>.
- Collins, M. and Brooks, J. (1995). Prepositional phrase attachment through a backed-off model. In *Workshop on Very Large Corpora*.
- Collins, M. and Duffy, N. (2001). Convolution kernels for natural language. In *Advances in neural information processing systems*, pages 625–632.
- Collins, M. and Singer, Y. (1999). Unsupervised models for named entity classification. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 189–196.
- Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 160–167.
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). *Introduction to algorithms*. MIT press, third edition.
- Coviello, L., Sohn, Y., Kramer, A. D., Marlow, C., Franceschetti, M., Christakis, N. A., and Fowler, J. H. (2014). Detecting emotional contagion in massive social networks. *PloS one*, 9(3):e90315.
- Crammer, K., Dekel, O., Keshet, J., Shalev-Shwartz, S., and Singer, Y. (2006). On-line passive-aggressive algorithms. *The Journal of Machine Learning Research*, 7:551–585.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Crammer, K. and Singer, Y. (2001). Pranking with ranking. In *Neural Information Processing Systems (NIPS)*, pages 641–647, Vancouver.
- Crammer, K. and Singer, Y. (2003). Ultraconservative online algorithms for multi-class problems. *The Journal of Machine Learning Research*, 3:951–991.
- Cui, H., Sun, R., Li, K., Kan, M.-Y., and Chua, T.-S. (2005). Question answering passage retrieval using dependency relations. In *Proceedings of the 28th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 400–407. ACM.
- De Marneffe, M.-C. and Manning, C. D. (2008). The stanford typed dependencies representation. In *Coling 2008: Proceedings of the workshop on Cross-Framework and Cross-Domain Parser Evaluation*, pages 1–8. Association for Computational Linguistics.
- Deerwester, S. C., Dumais, S. T., Landauer, T. K., Furnas, G. W., and Harshman, R. A. (1990). Indexing by latent semantic analysis. *JASIS*, 41(6):391–407.
- Derrida, J. et al. (1985). Des tours de babel. *Difference in translation*, 165.
- Duchi, J., Hazan, E., and Singer, Y. (2011). Adaptive subgradient methods for online learning and stochastic optimization. *The Journal of Machine Learning Research*, 12:2121–2159.
- Durrett, G. and Klein, D. (2013). Easy victories and uphill battles in coreference resolution. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*.
- Dyer, C. (2014). Notes on adagrad. www.ark.cs.cmu.edu/cdyer/adagrad.pdf.
- Eisner, J. and Satta, G. (1999). Efficient parsing for bilexical context-free grammars and head automaton grammars. In *Proceedings of the 37th annual meeting of the Association for Computational Linguistics on Computational Linguistics*, pages 457–464. Association for Computational Linguistics.
- Figueiredo, M., Graça, J., Martins, A., Almeida, M., and Coelho, L. P. (2013). LXMLS lab guide. <http://lxmls.it.pt/2013/guide.pdf>.
- Fillmore, C. J. (1968). The case for case. In *Universals in linguistic theory*. Holt, Rinehart, and Winston.

- Finkel, J. R., Grenager, T., and Manning, C. (2005). Incorporating non-local information into information extraction systems by gibbs sampling. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 363–370.
- Finkel, J. R., Kleeman, A., and Manning, C. D. (2008). Efficient, feature-based, conditional random field parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 959–967, Columbus, OH.
- Freund, Y., Schapire, R., and Abe, N. (1999). A short introduction to boosting. *Journal-Japanese Society For Artificial Intelligence*, 14(771-780):1612.
- Freund, Y. and Schapire, R. E. (1999). Large margin classification using the perceptron algorithm. *Machine learning*, 37(3):277–296.
- Gao, J., Andrew, G., Johnson, M., and Toutanova, K. (2007). A comparative study of parameter estimation methods for statistical natural language processing. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 824–831, Prague.
- Ge, D., Jiang, X., and Ye, Y. (2011). A note on the complexity of l_p minimization. *Mathematical programming*, 129(2):285–299.
- Goldwater, S. and Griffiths, T. (2007). A fully bayesian approach to unsupervised part-of-speech tagging. In *Annual meeting-association for computational linguistics*, volume 45.
- Grosz, B. J., Weinstein, S., and Joshi, A. K. (1995). Centering: A framework for modeling the local coherence of discourse. *Computational linguistics*, 21(2):203–225.
- Gruber, J. S. (1965). *Studies in Lexical Relations*. PhD thesis, Massachusetts Institute of Technology.
- Gundel, J. K., Hedberg, N., and Zacharski, R. (1993). Cognitive status and the form of referring expressions in discourse. *Language*, pages 274–307.
- Gutmann, M. U. and Hyvärinen, A. (2012). Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *The Journal of Machine Learning Research*, 13(1):307–361.
- Haghighi, A. and Klein, D. (2006). Prototype-driven learning for sequence models. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 320–327, New York, NY.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Hannak, A., Anderson, E., Barrett, L. F., Lehmann, S., Mislove, A., and Riedewald, M. (2012). Tweetin' in the rain: Exploring societal-scale effects of weather on mood. In *Proceedings of the International Conference on Web and Social Media (ICWSM)*, Menlo Park, California. AAAI Publications.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning*. Springer, New York, second edition.
- Hatzivassiloglou, V. and McKeown, K. R. (1997). Predicting the semantic orientation of adjectives. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 174–181, Madrid, Spain.
- Hindle, D. and Rooth, M. (1990). Structural ambiguity and lexical relations. In *Proceedings of the Workshop on Speech and Natural Language*.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Hsu, D., Kakade, S. M., and Zhang, T. (2012). A spectral algorithm for learning hidden markov models. *Journal of Computer and System Sciences*, 78(5):1460–1480.
- Hu, M. and Liu, B. (2004). Mining and summarizing customer reviews. In *Proceedings of Knowledge Discovery and Data Mining (KDD)*, pages 168–177.
- Huang, L., Fayong, S., and Guo, Y. (2012). Structured perceptron with inexact search. In *Proceedings of the 2012 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 142–151, Montréal, Canada. Association for Computational Linguistics.
- Ide, N. and Wilks, Y. (2006). Making sense about sense. In *Word sense disambiguation*, pages 47–73. Springer.
- Jaeger, T. and Levy, R. P. (2006). Speakers optimize information density through syntactic reduction. In *Neural Information Processing Systems (NIPS)*, pages 849–856.
- Jaeger, T. F. (2010). Redundancy and reduction: Speakers manage syntactic information density. *Cognitive psychology*, 61(1):23–62.
- Jain, A. K. (2010). Data clustering: 50 years beyond k-means. *Pattern recognition letters*, 31(8):651–666.

- Jockers, M. L. (2015). Szuzhet? <http://bla.bla.com>.
- Johnson, M. (1998). Pcfg models of linguistic tree representations. *Computational Linguistics*, 24(4):613–632.
- Jurafsky, D. and Martin, J. H. (2009). *Speech and Language Processing (2nd Edition) (Prentice Hall Series in Artificial Intelligence)*. Prentice Hall, 2 edition.
- Karlsson, F. (2007). Constraints on multiple center-embedding of clauses. *Journal of Linguistics*, 43(02):365–392.
- Khudanpur, S. and Wu, J. (2000). Maximum entropy techniques for exploiting syntactic, semantic and collocational dependencies in language modeling. *Computer Speech & Language*, 14(4):355–372.
- Kilgarrriff, A. (1997). I don’t believe in word senses. *CoRR*, cmp-lg/9712006.
- Kilgarrriff, A. and Grefenstette, G. (2003). Introduction to the special issue on the web as corpus. *Computational linguistics*, 29(3):333–347.
- Klein, D. and Manning, C. D. (2003). Accurate unlexicalized parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 423–430.
- Knight, K. and May, J. (2009). Applications of weighted automata in natural language processing. In *Handbook of Weighted Automata*, pages 571–596. Springer.
- Koo, T., Carreras, X., and Collins, M. (2008). Simple semi-supervised dependency parsing. In *Proceedings of ACL-08: HLT*, pages 595–603, Columbus, Ohio. Association for Computational Linguistics.
- Koo, T. and Collins, M. (2010). Efficient third-order dependency parsers. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1–11, Uppsala, Sweden.
- Koo, T., Globerson, A., Carreras, X., and Collins, M. (2007). Structured prediction models via the matrix-tree theorem. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 141–150.
- Koren, Y., Bell, R., and Volinsky, C. (2009). Matrix factorization techniques for recommender systems. *Computer*, 42(8):30–37.
- Lafferty, J., McCallum, A., and Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the International Conference on Machine Learning (ICML)*.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Lesk, M. (1986). Automatic sense disambiguation using machine readable dictionaries: how to tell a pine cone from an ice cream cone. In *Proceedings of the 5th annual international conference on Systems documentation*, pages 24–26. ACM.
- Levy, O. and Goldberg, Y. (2014). Neural word embedding as implicit matrix factorization. In *NIPS*.
- Levy, R. and Manning, C. (2009). An informal introduction to computational semantics.
- Liang, P. and Klein, D. (2009). Online em for unsupervised models. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 611–619, Boulder, CO.
- Ling, W., Luís, T., Marujo, L., Astudillo, R. F., Amir, S., Dyer, C., Black, A. W., and Trancoso, I. (2015). Finding function in form: Compositional character models for open vocabulary word representation.
- Liu, B. (2012). Sentiment analysis and opinion mining. *Synthesis Lectures on Human Language Technologies*, 5(1):1–167.
- Liu, B. (2015). *Sentiment Analysis: Mining Opinions, Sentiments, and Emotions*. Cambridge University Press.
- Liu, D. C. and Nocedal, J. (1989). On the limited memory BFGS method for large scale optimization. *Mathematical programming*, 45(1-3):503–528.
- Manning, C. D., Raghavan, P., Schütze, H., et al. (2008). *Introduction to information retrieval*, volume 1. Cambridge university press Cambridge.
- Manning, C. D. and Schütze, H. (1999). *Foundations of Statistical Natural Language Processing*. MIT press, Cambridge, Massachusetts.
- Marcus, M. P., Marcinkiewicz, M. A., and Santorini, B. (1993). Building a large annotated corpus of English: The Penn Treebank. *Computational Linguistics*, 19(2):313–330.
- McClosky, D., Charniak, E., and Johnson, M. (2006). Effective self-training for parsing. In *Proceedings of the main conference on human language technology conference of the North American Chapter of the Association of Computational Linguistics*, pages 152–159. Association for Computational Linguistics.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- McLendon, S. (2003). Evidentials in eastern pomo with a comparative survey of the category in other pomoan languages. *TYPOLOGICAL STUDIES IN LANGUAGE*, 54:101–130.
- Mihalcea, R., Chklovski, T. A., and Kilgarrieff, A. (2004). The senseval-3 english lexical sample task. In *Proceedings of SENSEVAL-3*, pages 25–28, Barcelona, Spain. Association for Computational Linguistics.
- Mikolov, T., Chen, K., Corrado, G., and Dean, J. (2013a). Efficient estimation of word representations in vector space. In *Proceedings of International Conference on Learning Representations*.
- Mikolov, T., Deoras, A., Povey, D., Burget, L., and Cernocky, J. (2011). Strategies for training large scale neural network language models. In *Automatic Speech Recognition and Understanding (ASRU), 2011 IEEE Workshop on*, pages 196–201. IEEE.
- Mikolov, T., Karafiát, M., Burget, L., Cernocký, J., and Khudanpur, S. (2010). Recurrent neural network based language model. In *Proceedings of the International Speech Communication Association (INTERSPEECH)*, pages 1045–1048.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013b). Distributed representations of words and phrases and their compositionality. In *Advances in Neural Information Processing Systems*, pages 3111–3119.
- Mikolov, T., Yih, W.-t., and Zweig, G. (2013c). Linguistic regularities in continuous space word representations. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 746–751, Stroudsburg, Pennsylvania. Association for Computational Linguistics.
- Miller, M., Sathi, C., Wiesensthal, D., Leskovec, J., and Potts, C. (2011). Sentiment flow through hyperlink networks. In *Proceedings of the International Conference on Web and Social Media (ICWSM)*, Menlo Park, California. AAAI Publications.
- Miller, S., Guinness, J., and Zamanian, A. (2004). Name tagging with word clusters and discriminative training. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 337–342, Boston, MA.
- Minka, T. P. (1999). From hidden markov models to linear dynamical systems. Tech. Rep. 531, Vision and Modeling Group of Media Lab, MIT.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Mitra, S. and Acharya, T. (2007). Gesture recognition: A survey. *Systems, Man, and Cybernetics, Part C: Applications and Reviews, IEEE Transactions on*, 37(3):311–324.
- Mnih, A. and Hinton, G. E. (2008). A scalable hierarchical distributed language model. In *Neural Information Processing Systems (NIPS)*, pages 1081–1088.
- Mnih, A. and Teh, Y. W. (2012). A fast and simple algorithm for training neural probabilistic language models. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Mohri, M., Pereira, F., and Riley, M. (2002). Weighted finite-state transducers in speech recognition. *Computer Speech & Language*, 16(1):69–88.
- Mohri, M., Rostamizadeh, A., and Talwalkar, A. (2012). *Foundations of machine learning*. MIT press.
- Muralidharan, A. and Hearst, M. A. (2013). Supporting exploratory text analysis in literature study. *Literary and linguistic computing*, 28(2):283–295.
- Murphy, K. P. (2012). *Machine Learning: A Probabilistic Perspective*. The MIT Press.
- Nakagawa, T., Inui, K., and Kurohashi, S. (2010). Dependency tree-based sentiment classification using crfs with hidden variables. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 786–794, Los Angeles, CA.
- Navigli, R. (2009). Word sense disambiguation: A survey. *ACM Computing Surveys (CSUR)*, 41(2):10.
- Neal, R. M. and Hinton, G. E. (1998). A view of the em algorithm that justifies incremental, sparse, and other variants. In *Learning in graphical models*, pages 355–368. Springer.
- Nemirovski, A. and Yudin, D. (1978). On Cezari’s convergence of the steepest descent method for approximating saddle points of convex-concave functions. *Soviet Math. Dokl.*, 19.
- Neuhaus, P. and Bröker, N. (1997). The complexity of recognition of linguistically adequate dependency grammars. In *Proceedings of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 337–343.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Nigam, K., McCallum, A. K., Thrun, S., and Mitchell, T. (2000). Text classification from labeled and unlabeled documents using em. *Machine learning*, 39(2-3):103–134.
- Nivre, J. (2004). Incrementality in deterministic dependency parsing. In *Proceedings of the Workshop on Incremental Parsing: Bringing Engineering and Cognition Together*, pages 50–57. Association for Computational Linguistics.
- Nivre, J. and Nilsson, J. (2005). Pseudo-projective dependency parsing. In *Proceedings of the 43rd Annual Meeting on Association for Computational Linguistics*, pages 99–106. Association for Computational Linguistics.
- Och, F. J. and Ney, H. (2003). A systematic comparison of various statistical alignment models. *Computational linguistics*, 29(1):19–51.
- Pak, A. and Paroubek, P. (2010). Twitter as a corpus for sentiment analysis and opinion mining. In *Proceedings of the Language Resources and Evaluation Conference*, volume 10, pages 1320–1326.
- Palmer, M., Gildea, D., and Kingsbury, P. (2005). The proposition bank: An annotated corpus of semantic roles. *Computational linguistics*, 31(1):71–106.
- Pang, B. and Lee, L. (2005). Seeing stars: Exploiting class relationships for sentiment categorization with respect to rating scales. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 115–124.
- Pang, B. and Lee, L. (2008). Opinion mining and sentiment analysis. *Foundations and trends in information retrieval*, 2(1-2):1–135.
- Pang, B., Lee, L., and Vaithyanathan, S. (2002). Thumbs up?: sentiment classification using machine learning techniques. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 79–86.
- Pantel, P. and Lin, D. (2002). Discovering word senses from text. In *Proceedings of Knowledge Discovery and Data Mining (KDD)*, pages 613–619. ACM.
- Pedersen, T. and Bruce, R. (1997). Distinguishing word senses in untagged text. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global vectors for word representation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1532–1543.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Pereira, F. (2000). Formal grammar and information theory: together again? *Philosophical Transactions of the Royal Society of London. Series A: Mathematical, Physical and Engineering Sciences*, 358(1769):1239–1253.
- Pereira, F., Tishby, N., and Lee, L. (1993). Distributional clustering of english words. In *Proceedings of the 31st annual meeting on Association for Computational Linguistics*, pages 183–190. Association for Computational Linguistics.
- Petrov, S., Barrett, L., Thibaux, R., and Klein, D. (2006). Learning accurate, compact, and interpretable tree annotation. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Petrov, S., Das, D., and McDonald, R. (2012). A universal part-of-speech tagset. In *Proceedings of LREC*.
- Petrov, S. and Klein, D. (2007). Improved inference for unlexicalized parsing. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 404–411.
- Rabiner, L. (1989). A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286.
- Raghunathan, K., Lee, H., Rangarajan, S., Chambers, N., Surdeanu, M., Jurafsky, D., and Manning, C. (2010). A multi-pass sieve for coreference resolution. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 492–501, Stroudsburg, Pennsylvania. Association for Computational Linguistics.
- Rao, D. and Ravichandran, D. (2009). Semi-supervised polarity lexicon induction. In *Proceedings of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 675–682.
- Ratnaparkhi, A. (1999). Learning to parse natural language with maximum entropy models. *Machine learning*, 34(1-3):151–175.
- Ratnaparkhi, A., Reynar, J., and Roukos, S. (1994). A maximum entropy model for prepositional phrase attachment. In *Proceedings of the workshop on Human Language Technology*, pages 250–255. Association for Computational Linguistics.
- Read, J. (2005). Using emoticons to reduce dependency in machine learning techniques for sentiment classification. In *Proceedings of the ACL student research workshop*, pages 43–48. Association for Computational Linguistics.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Reisinger, J. and Mooney, R. J. (2010). Multi-prototype vector-space models of word meaning. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 109–117, Los Angeles, CA.
- Resnik, P. and Smith, N. A. (2003). The web as a parallel corpus. *Computational Linguistics*, 29(3):349–380.
- Riloff, E. (1996). Automatically generating extraction patterns from untagged text. In *Proceedings of the national conference on artificial intelligence*, pages 1044–1049.
- Roark, B., Saraclar, M., and Collins, M. (2007). Discriminative n -gram language modeling. *Computer Speech & Language*, 21(2):373–392.
- Robert, C. and Casella, G. (2013). *Monte Carlo statistical methods*. Springer Science & Business Media.
- Rosenfeld, R. (1996). A maximum entropy approach to adaptive statistical language modelling. *Computer Speech & Language*, 10(3):187–228.
- Samdani, R., Chang, M.-W., and Roth, D. (2012). Unified expectation maximization. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 688–698.
- Sato, M.-A. and Ishii, S. (2000). On-line em algorithm for the normalized gaussian network. *Neural computation*, 12(2):407–432.
- Schmid, H. (1994). Probabilistic part-of-speech tagging using decision trees. In *Proceedings of the international conference on new methods in language processing*, volume 12, pages 44–49. Manchester, UK.
- Shalev-Shwartz, S., Singer, Y., and Srebro, N. (2007). Pegasos: Primal estimated sub-GrAdient SOLver for SVM. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 807–814.
- Shen, D. and Lapata, M. (2007). Using semantic roles to improve question answering. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 12–21.
- Shieber, S. M. (1985). Evidence against the context-freeness of natural language. *Linguistics and Philosophy*, 8(3):333–343.
- Sipser, M. (2012). *Introduction to the Theory of Computation*. Cengage Learning.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Smith, N. A. (2011). Linguistic structure prediction. *Synthesis Lectures on Human Language Technologies*, 4(2):1–274.
- Smith, N. J. and Levy, R. (2013). The effect of word predictability on reading time is logarithmic. *Cognition*, 128(3):302–319.
- Socher, R., Bauer, J., Manning, C. D., and Ng, A. Y. (2013). Parsing With Compositional Vector Grammars. In *Proceedings of the Association for Computational Linguistics (ACL)*, Sophia, Bulgaria.
- Song, L., Boots, B., Siddiqi, S. M., Gordon, G. J., and Smola, A. J. (2010). Hilbert space embeddings of hidden markov models. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 991–998.
- Soon, W. M., Ng, H. T., and Lim, D. C. Y. (2001). A machine learning approach to coreference resolution of noun phrases. *Computational linguistics*, 27(4):521–544.
- Spitkovsky, V. I., Alshawih, H., Jurafsky, D., and Manning, C. D. (2010). Viterbi training improves unsupervised dependency parsing. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 9–17.
- Sproat, R., Black, A., Chen, S., Kumar, S., Ostendorf, M., and Richards, C. (2001). Normalization of non-standard words. *Computer Speech & Language*, 15(3):287–333.
- Sra, S., Nowozin, S., and Wright, S. J. (2012). *Optimization for machine learning*. MIT Press.
- Sundermeyer, M., Schlüter, R., and Ney, H. (2012). Lstm neural networks for language modeling. In *Proceedings of the International Speech Communication Association (INTERSPEECH)*.
- Taskar, B., Guestrin, C., and Koller, D. (2003). Max-margin markov networks. In *NIPS*.
- Tausczik, Y. R. and Pennebaker, J. W. (2010). The psychological meaning of words: LIWC and computerized text analysis methods. *Journal of Language and Social Psychology*, 29(1):24–54.
- Teh, Y. W. (2006). A hierarchical bayesian language model based on pitman-yor processes. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 985–992.

- Thomas, M., Pang, B., and Lee, L. (2006). Get out the vote: Determining support or opposition from Congressional floor-debate transcripts. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 327–335.
- Tibshirani, R. (1996). Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 267–288.
- Tsochantaridis, I., Hofmann, T., Joachims, T., and Altun, Y. (2004). Support vector machine learning for interdependent and structured output spaces. In *Proceedings of the twenty-first international conference on Machine learning*, page 104. ACM.
- Turney, P. (2003). Thumbs up or thumbs down? semantic orientation applied to unsupervised classification of reviews. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 417–424.
- Van Gael, J., Vlachos, A., and Ghahramani, Z. (2009). The infinite hmm for unsupervised pos tagging. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 678–687, Singapore.
- Wei, G. C. and Tanner, M. A. (1990). A Monte Carlo implementation of the EM algorithm and the poor man’s data augmentation algorithms. *Journal of the American Statistical Association*, 85(411):699–704.
- Wilson, T., Wiebe, J., and Hoffmann, P. (2005). Recognizing contextual polarity in phrase-level sentiment analysis. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 347–354.
- Xu, W., Liu, X., and Gong, Y. (2003). Document clustering based on non-negative matrix factorization. In *Proceedings of ACM SIGIR conference on Research and development in information retrieval*, pages 267–273. ACM.
- Yamada, H. and Matsumoto, Y. (2003). Statistical dependency analysis with support vector machines. In *Proceedings of IWPT*, volume 3, pages 195–206.
- Yarowsky, D. (1995). Unsupervised word sense disambiguation rivaling supervised methods. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 189–196. Association for Computational Linguistics.
- Zaidan, O. F. and Callison-Burch, C. (2011). Crowdsourcing translation: Professional quality from non-professionals. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1220–1229, Portland, OR.

(c) Jacob Eisenstein 2014-2015. Work in progress.

- Zelenko, D., Aone, C., and Richardella, A. (2003). Kernel methods for relation extraction. *The Journal of Machine Learning Research*, 3:1083–1106.
- Zhang, T. (2004). Solving large scale linear prediction problems using stochastic gradient descent algorithms. In *Proceedings of the twenty-first international conference on Machine learning*, page 116. ACM.
- Zhu, X. and Ghahramani, Z. (2002). Learning from labeled and unlabeled data with label propagation. Technical report, Technical Report CMU-CALD-02-107, Carnegie Mellon University.