

Natural Language Processing

Jacob Eisenstein

December 21, 2017

Contents

Contents	1
I Words, bags of words, and features	9
1 Linear text classification	11
1.1 Naïve Bayes	14
1.2 Discriminative learning	21
1.3 Loss functions and large-margin classification	25
1.4 Logistic regression	31
1.5 Optimization	33
1.6 *Additional topics in classification	36
1.7 Summary of learning algorithms	38
2 Nonlinear classification	41
2.1 Feedforward neural networks	42
2.2 Neural networks	44
2.3 Learning neural networks	47
2.4 Convolutional neural networks	55
3 Linguistic applications of classification	61
3.1 Sentiment and opinion analysis	61
3.2 Word sense disambiguation	65
3.3 Design decisions for text classification	68
3.4 Evaluating classifiers	72
3.5 Building datasets	79
4 Learning without supervision	87
4.1 Unsupervised learning	87
4.2 Applications of expectation-maximization	96
4.3 Semi-supervised learning	99

4.4	Domain adaptation	102
4.5	*Other approaches to learning with latent variables	106
II	Sequences and trees	113
5	Language models	115
5.1	N-gram language models	117
5.2	Smoothing and discounting	119
5.3	Recurrent neural network language models	124
5.4	Evaluating language models	130
5.5	Out-of-vocabulary words	132
6	Sequence labeling	135
6.1	Sequence labeling as classification	135
6.2	Sequence labeling as structure prediction	137
6.3	The Viterbi algorithm	139
6.4	Hidden Markov Models	143
6.5	Discriminative sequence labeling	150
6.6	*Unsupervised sequence labeling	162
7	Applications of sequence labeling	167
7.1	Part-of-speech tagging	167
7.2	Morphosyntactic Attributes	173
7.3	Named entity recognition	175
7.4	Tokenization	177
7.5	Code switching	177
7.6	Dialogue acts	178
8	Formal language theory	181
8.1	Regular languages	182
8.2	Context-free languages	197
8.3	*Mildly context-sensitive languages	210
9	Context-free Parsing	217
9.1	Deterministic bottom-up parsing	217
9.2	Ambiguity in parsing	220
9.3	Weighted Context-Free Grammars	224
9.4	Improving Parsing by Refined Non-terminals	229
9.5	Discriminative parsing	238
10	Dependency Parsing	241

10.1	Dependency grammar	241
10.2	Graph-based dependency parsing	246
10.3	Transition-based dependency parsing	253
10.4	Applications	257
III Meaning		259
11	Compositional semantics	261
11.1	Meaning and denotation	262
11.2	Logical representations of meaning	262
11.3	Semantic parsing and the lambda calculus	267
11.4	Learning semantic parsers	272
12	Predicate-argument semantics	283
12.1	Semantic roles	285
12.2	Semantic role labeling	290
12.3	Abstract Meaning Representation	299
12.4	Applications of Predicate-Argument Semantics	303
13	Distributional and distributed semantics	309
13.1	The distributional hypothesis	309
13.2	Design decisions for word representations	311
13.3	Latent semantic analysis	314
13.4	Brown clusters	315
13.5	Neural word embeddings	318
13.6	Evaluating word embeddings	323
13.7	Distributed representations beyond distributional statistics	325
13.8	Distributed representations of multiword units	328
14	Reference Resolution	335
14.1	Forms of referring expressions	336
14.2	Learning for coreference resolution	343
14.3	Cross-document coreference resolution	347
15	Discourse	349
15.1	Discourse relations in the Penn Discourse Treebank	349
15.2	Rhetorical Structure Theory	349
15.3	Centering	349

IV Applications	351
16 Information extraction	353
16.1 Entities	355
16.2 Relations	361
16.3 Events	370
16.4 Hedges, denials, and hypotheticals	372
16.5 Question answering and machine reading	374
17 Text generation	381
18 Machine translation	383
A Probability	385
A.1 Probabilities of event combinations	386
A.2 Conditional probability and Bayes' rule	387
A.3 Independence	389
A.4 Random variables	390
A.5 Expectations	391
A.6 Modeling and estimation	392
B Computational complexity	395
Bibliography	397

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. There are 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 time 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 machine learning presentations, due to the 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.

The notes focus 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. A remarkable thing about the field of natural language processing is that so many problems in language technology can be solved by a small number of methods. These notes focus on the following methods:

Search algorithms shortest path, Viterbi, CKY, minimum spanning tree, shift-reduce, integer linear programming, beam search.

Learning algorithms Naïve Bayes, logistic regression, perceptron, expectation-maximization, matrix factorization, backpropagation.

The goal of this text is to teach how these methods work, and how they can be applied to problems that arise in the computer processing of natural language: 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. Proper application of these techniques requires understanding the underlying linguistic phenomena, and the notes therefore include a minimal foundation in morphology, syntax, semantics, and discourse. However, a detailed understanding of these topics can be provided only by a linguistics textbook (e.g., Akmajian et al., 2010; Fromkin et al., 2013).

-Jacob Eisenstein, December 21, 2017

Notation

As a general rule, random variables and observable counts are indicated with Roman letters (a, b, c), and parameters are indicated with Greek letters (α, β, θ). Vectors are indicated with bold script for both random variables \mathbf{x} and parameters $\boldsymbol{\theta}$. Other useful notations are indicated in the table below.

Linear algebra	
$\mathbf{x}^{(i)}$	a column vector of feature counts for instance i , often word counts
$\mathbf{x}_{j:k}$	elements j through k (inclusive) of a vector \mathbf{x}
$[\mathbf{x}; \mathbf{y}]$	vertical concatenation of two column vectors
$[\mathbf{x}, \mathbf{y}]$	horizontal concatenation of two column vectors
\mathbf{e}_n	a “one-hot” vector with a value of 1 at position n , and zero everywhere else
$\boldsymbol{\theta}^\top$	the transpose of a column vector $\boldsymbol{\theta}$
$\boldsymbol{\theta} \cdot \mathbf{x}^{(i)}$	the dot product of vectors $\boldsymbol{\theta}$ and $\mathbf{x}^{(i)}$
\mathbf{X}	a matrix
$x_{i,j}$	row i , column j of matrix \mathbf{X}
Text datasets	
w_m	word token at position m
N	number of training instances
M	length of a sequence (of words or tags)
V	number of words in vocabulary
$y^{(i)}$	the true label for instance i
\hat{y}	a predicted label
\mathcal{Y}	the set of all possible labels
K	number of possible labels $K = \mathcal{Y} $

Probabilities

$\Pr(A)$	probability of event A
$\Pr(A \mid B)$	probability of event A , conditioned on event B
$p_B(b)$	the marginal probability of random variable B taking value b ; written $p(b)$ when the choice of random variable is clear from context
$p_{B A}(b \mid a)$	the probability of random variable B taking value b , conditioned on A taking value a ; written $p(b \mid a)$ when clear from context
$a \mid b \sim f(\theta_b)$	conditioned on b , a is drawn from distribution of form f with parameters θ_b [todo: decide about upper-case random variables]

Linear classification

$\mathbf{f}(\mathbf{x}^{(i)}, y)$	feature vector for instance i with label y
$\boldsymbol{\theta}$	a (column) vector of weights
$\ell^{(i)}$	loss on an individual instance i
L	objective function for an entire dataset
\mathcal{L}	log-likelihood of a dataset
λ	the amount of regularization

Sequence labeling

$\mathcal{Y}(\mathbf{w})$	the set of possible tag sequences for the word sequence \mathbf{w}
\diamond	the start tag
\blacklozenge	the stop tag
\square	the start token
\blacksquare	the stop token

Part I

Words, bags of words, and features

Chapter 1

Linear text classification

We'll start with the problem of **text classification**: given a text document, assign it a discrete label $y \in \mathcal{Y}$, where \mathcal{Y} is the set of possible labels. This problem has many applications, from spam filtering to analysis of electronic health records. Text classification is also a building block that is used throughout more complex natural language processing tasks.

To perform this task, the first question is how to represent each document. A common approach is to use a vector of word counts, e.g., $\mathbf{x}^{(i)} = [0 \ 1 \ 1 \ 0 \ 0 \ 2 \ 0 \ 1 \ 13 \ 0 \ \dots]^\top$, where $x_j^{(i)}$ is the count of word j in document i . The length of $\mathbf{x}^{(i)}$ is $V \triangleq |\mathcal{V}|$, where \mathcal{V} is the set of possible words in the vocabulary.

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 the words appear. We have thrown out grammar, sentence boundaries, paragraphs — everything but the words. Yet the bag of words model is surprisingly effective for text classification. If you see the word *freeee* in an email, is it a spam email? What if you see the word *Bayesian*? For many labeling problems, individual words can be strong predictors.

To predict a label from a bag-of-words, we can assign a score to each word in the vocabulary, measuring the compatibility with the label. In the spam filtering case, we might assign a positive score to the word *freeee* for the label SPAM, and a negative score to the word *Bayesian*. These scores are called **weights**, written $\boldsymbol{\theta}$.

Now suppose we want a multiclass classifier, where $K \triangleq |\mathcal{Y}| > 2$. For example, we might want to classify news stories about sports, celebrities, music, and business. The goal is to predict a label $\hat{y}^{(i)}$, given the bag of words $\mathbf{x}^{(i)}$, using the weights $\boldsymbol{\theta}$. We will do this using a vector inner product between the weights $\boldsymbol{\theta}$ and a **feature function** $\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$. As the notation suggests, the feature vector is a *function* of two arguments, the word counts $\mathbf{x}^{(i)}$ and the label $y^{(i)}$, and it returns a vector output. For example, for arguments $\mathbf{x}^{(i)}$ and

y , element j of this feature vector might be,

$$f_j(\mathbf{x}, y) = \begin{cases} x_{\text{freee}}, & \text{if } y = \text{SPAM} \\ 0, & \text{otherwise} \end{cases} \quad [1.1]$$

This function returns the count of the word *freee* if the label is SPAM, and it returns zero otherwise. The corresponding weight θ_j then scores the compatibility of the word *freee* with the label SPAM. A positive score means that this word makes the label more likely.

Let's now formalize this feature function. For any pair $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$, we define $\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$ as,

$$\mathbf{f}(\mathbf{x}, y = 0) = [\mathbf{x}; \underbrace{0; 0; \dots; 0}_{(K-1) \times V}] \quad [1.2]$$

$$\mathbf{f}(\mathbf{x}, y = 1) = [\underbrace{0; 0; \dots; 0}_V; \mathbf{x}; \underbrace{0; 0; \dots; 0}_{(K-2) \times V}] \quad [1.3]$$

$$\mathbf{f}(\mathbf{x}, y = 2) = [\underbrace{0; 0; \dots; 0}_{2 \times V}; \mathbf{x}; \underbrace{0; 0; \dots; 0}_{(K-3) \times V}] \quad [1.4]$$

$$\mathbf{f}(\mathbf{x}, y = K) = [\underbrace{0; 0; \dots; 0}_{(K-1) \times V}; \mathbf{x}], \quad [1.5]$$

where $\underbrace{[0; 0; \dots; 0]}_{(K-1) \times V}$ is a column vector of $(K-1) \times V$ zeros, and the semicolon indicates vertical concatenation. This arrangement is shown in Figure 1.1; the notation may seem awkward at first, but it helps to keep things simple as we progress to more complex learning settings.

Given a vector of weights, $\boldsymbol{\theta} \in \mathbb{R}^{V \times K}$, we can now compute the inner product $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y)$. This inner product gives a scalar measure of the compatibility of the observation \mathbf{x} with label y .¹ For any document $\mathbf{x}^{(i)}$, we predict the label \hat{y} ,

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) \quad [1.6]$$

This inner product notation gives a clean separation between the *data* (\mathbf{x} and y) and the *parameters* ($\boldsymbol{\theta}$). As we will see in later chapters, this notation also generalizes nicely to **structured classification**, in which the space of labels \mathcal{Y} is very large, and we want to model shared substructures between labels.

¹Only $V \times (K-1)$ features and weights are necessary: by stipulating that $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y = K) = 0$ regardless of \mathbf{x} , it is possible to implement any classification rule that can be achieved with $V \times K$ features and weights. This is the approach taken in binary classification rules like $y = \text{Sign}(\boldsymbol{\beta} \cdot \mathbf{x} + a)$, where $\boldsymbol{\beta}$ is a vector of weights, a is an offset, and the label set is $\mathcal{Y} = \{-1, 1\}$. However, for multiclass classification, it is more concise to write $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y)$ for all $y \in \mathcal{Y}$.

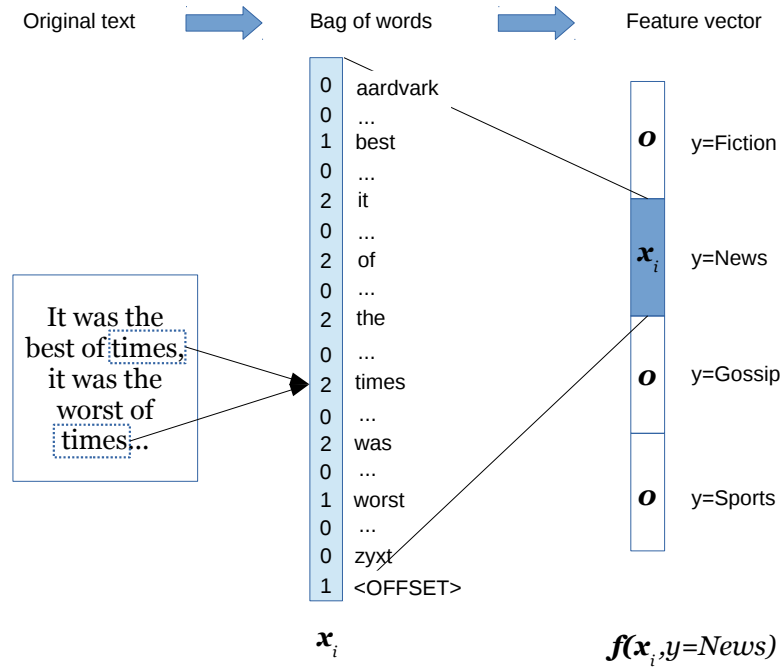


Figure 1.1: The bag-of-words and feature vector representations, for a hypothetical text classification task.

It is common to add an **offset feature** at the end of the vector of word counts x , which is always 1. We then have to also add an extra zero to each of the zero vectors, to make the vector lengths match. This gives the entire feature vector $f(x, y)$ a length of $(V + 1) \times K$. The weight associated with this offset feature can be thought of as a bias for (or against) 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 θ , where do they come from? We could set them by hand. If we wanted to distinguish, say, English from Spanish, we can use English and Spanish dictionaries, and set the weight to one for each word that appears in the associated dictionary. For example,²

$$\begin{aligned}
 \theta_{(E, \text{bicycle})} &= 1 & \theta_{(S, \text{bicycle})} &= 0 \\
 \theta_{(E, \text{bicicleta})} &= 0 & \theta_{(S, \text{bicicleta})} &= 1 \\
 \theta_{(E, \text{con})} &= 1 & \theta_{(S, \text{con})} &= 1 \\
 \theta_{(E, \text{ordinateur})} &= 0 & \theta_{(S, \text{ordinateur})} &= 0.
 \end{aligned}$$

²In this notation, each tuple (language, word) indexes an element in θ , which remains a vector.

Similarly, if we want to distinguish positive and negative sentiment, we could use positive and negative **sentiment lexicons**, which are defined by social psychologists (Tausczik and Pennebaker, 2010).

But it is usually not easy to set classification weights by hand. Instead, we will learn them from data. Email users manually label messages as SPAM; newspapers label their own articles as BUSINESS or STYLE. Using such **instance labels**, we can automatically acquire weights using **supervised machine learning**. This chapter will discuss several machine learning approaches for classification. The first is based on probability. (For a review of probability, consult Appendix A.)

1.1 Naïve Bayes

The **joint probability** of a bag of words $\mathbf{x}^{(i)}$ and its true label $y^{(i)}$ is written $p(\mathbf{x}^{(i)}, y^{(i)})$. Suppose we have a dataset of N labeled instances, $\{\mathbf{x}^{(1:N)}, y^{(1:N)}\}$, which we assume are **independent and identically distributed (IID)** (see § A.3). Then the joint probability of the entire dataset is,

$$p(\mathbf{x}^{(1:N)}, y^{(1:N)}) = \prod_{i=1}^N p(\mathbf{x}^{(i)}, y^{(i)}). \quad [1.7]$$

What does this have to do with classification? One approach to classification is to set the weights θ so as to maximize the joint probability of a **training set** of labeled documents. This is known as **maximum likelihood estimation**:

$$\theta = \underset{\theta}{\operatorname{argmax}} p(\mathbf{x}^{(1:N)}, y^{(1:N)}; \theta) \quad [1.8]$$

$$= \underset{\theta}{\operatorname{argmax}} \prod_{i=1}^N p(\mathbf{x}^{(i)}, y^{(i)}; \theta) \quad [1.9]$$

$$= \underset{\theta}{\operatorname{argmax}} \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, y^{(i)}; \theta). \quad [1.10]$$

The notation $p(\mathbf{x}^{(i)}, y^{(i)}; \theta)$ indicates that θ is a *parameter* of the probability function. We can move from a product of probabilities to a sum of log probabilities because the log function is monotonically increasing over positive arguments, and so the same θ will maximize both the probability and the log probability. Working with logarithms is desirable because of numerical stability: on a large dataset, multiplying many probabilities can **underflow** to zero.

The probability $p(\mathbf{x}^{(i)}, y^{(i)}; \theta)$ is defined through a **generative model** — an idealized

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 1 Generative process for the Naïve Bayes classifier

for Document $i \in \{1, 2, \dots, N\}$ **do**:
 Draw the label $y^{(i)} \sim \text{Categorical}(\boldsymbol{\mu})$;
 Draw the word counts $\mathbf{x}^{(i)} \mid y^{(i)} \sim \text{Multinomial}(\boldsymbol{\phi}_{y^{(i)}})$.

random process that has generated the observed data, $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$.³ Algorithm 1 describes the generative model describes the **Naïve Bayes** classifier, with parameters $\boldsymbol{\theta} = \{\boldsymbol{\mu}, \boldsymbol{\phi}\}$.

- The first line of this generative model encodes the assumption that the instances are mutually independent: neither the label nor the text of document i affects the label or text of document j .⁴ Furthermore, the instances are identically distributed: the distributions over the label $y^{(i)}$ and the text $\mathbf{x}^{(i)}$ are the same for all instances i .
- The second line of the generative model states that the random variable $y^{(i)}$ is drawn from a categorical distribution with parameter $\boldsymbol{\mu}$. are just like weighted dice: the vector $\boldsymbol{\mu} = [\mu_1, \mu_2, \dots, \mu_K]$ gives the probabilities of each label, so that the probability of drawing label y is equal to μ_y . For example, if $\mathcal{Y} = \{\text{POSITIVE}, \text{NEGATIVE}, \text{NEUTRAL}\}$, we might have $\boldsymbol{\mu} = [0.1, 0.7, 0.2]$. We require $\sum_y \mu_y = 1$ and $\forall_y (\mu_y \geq 0)$.⁵
- The third line describes how the bag-of-words counts $\mathbf{x}^{(i)}$ are generated. By writing $\mathbf{x}^{(i)} \mid y^{(i)}$, this line indicates that the word counts are conditioned on the label, so that the joint probability is factored using the chain rule,

$$p(\mathbf{x}^{(i)}, y^{(i)}) = p(\mathbf{x}^{(i)} \mid y^{(i)}) \times p(y^{(i)}). \quad [1.11]$$

The specific distribution $p_{\mathbf{x}^{(i)} \mid y^{(i)}}$ is the **multinomial**, which is a probability distribution over vectors of non-negative counts. The probability mass function for this distribution is:

$$p_{\text{mult}}(\mathbf{x}; \boldsymbol{\phi}) = B(\mathbf{x}) \prod_{j=1}^V \phi_j^{x_j} \quad [1.12]$$

$$B(\mathbf{x}) = \frac{\left(\sum_{j=1}^V x_j\right)!}{\prod_{j=1}^V (x_j!)} \quad [1.13]$$

³Generative models will be used throughout this text. They explicitly define the assumptions underlying the form of a probability distribution over observed and latent variables. For a readable introduction to generative models in statistics, see Blei (2014).

⁴Can you think of any cases in which this assumption is too strong?

⁵Formally, we require $\boldsymbol{\mu} \in \Delta^{K-1}$, where Δ^{K-1} is the $K - 1$ **probability simplex**, the set of all vectors of K nonnegative numbers that sum to one. Because of the sum-to-one constraint, there are $K - 1$ degrees of freedom for a vector of size K .

Algorithm 2 Alternative generative process for the Naïve Bayes classifier

```

for Document  $i \in \{1, 2, \dots, N\}$  do:
  Draw the label  $y^{(i)} \sim \text{Categorical}(\mu)$ ;
  for Token  $m \in \{1, 2, \dots, M_i\}$  do:
    Draw the token  $w_m^{(i)} \mid y^{(i)} \sim \text{Categorical}(\phi_{y^{(i)}})$ .

```

As in the categorical distribution, the parameter ϕ_j can be interpreted as a probability: specifically, the probability that any given token in the document is the word j .⁶ The multinomial distribution involves a product over words, with each term in the product equal to the probability of the word ϕ_j , exponentiated by the count x_j . Words that have zero count play no role in this product, because $\phi_j^0 = 1$ for all ϕ_j . The term $B(\mathbf{x})$ doesn't depend on ϕ , and can usually be ignored. Can you see why we need this term at all?⁷

The notation $p(\mathbf{x} \mid y; \phi)$ indicates the conditional probability of word counts \mathbf{x} given label y , with parameter ϕ , which is equal to $p_{\text{mult}}(\mathbf{x}; \phi_y)$. By specifying the multinomial distribution, we describe the **multinomial naïve Bayes** classifier. Why “naïve”? Because the multinomial distribution treats each word token independently: the probability mass function factorizes across the counts.⁸

1.1.1 Types and tokens

A slight modification to the generative model of Naïve Bayes is shown in Algorithm 2. Instead of generating a vector of counts of **types**, \mathbf{x} , this model generates a *sequence of tokens*, $\mathbf{w} = (w_1, w_2, \dots, w_M)$. The distinction between types and tokens is critical: $x_j \in \{0, 1, 2, \dots, M\}$ is the count of word type j in the vocabulary, e.g., the number of times the word *cannibal* appears; $w_m \in \mathcal{V}$ is the identity of token m in the document, e.g. $w_m = \text{cannibal}$.

The probability of the sequence \mathbf{w} is a product of categorical probabilities. Algorithm 2 makes a conditional independence assumption: each token $w_m^{(i)}$ is independent of all other tokens $w_n^{(i)}$, conditioned on the label $y^{(i)}$. This is identical to the “naïve” independence assumption implied by the multinomial distribution, and as a result, the optimal param-

⁶You are encouraged to verify that the multinomial distribution is identical to the categorical distribution when $\sum_j x_j = 1$.

⁷Technically, a multinomial distribution requires a second parameter, the total number of word counts in \mathbf{x} . In the bag-of-words representation is equal to the number of words in the document. However, this parameter is irrelevant for classification.

⁸You 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 is naïve in exactly the same sense.

eters for this model are identical to those in multinomial Naïve Bayes. For any instance, the probability assigned by this model is proportional the probability under multinomial Naïve Bayes, except for the factor $B(\mathbf{x})$, which appears in the multinomial distribution. Because $B(\mathbf{x}) \geq 1$, 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: there can be many word sequences that correspond to a single vector of 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 $B(\mathbf{x})$ is equal to the total number of possible word orderings for count vector \mathbf{x} .

Sometimes it is useful to think of instances as counts of types, \mathbf{x} ; other times, it is better to think of them as sequences of tokens, \mathbf{w} . If the tokens are generated from a model that assumes independence, conditioned on the label, then these two views lead to probability models that are identical, except for a scaling factor that does not depend on the label or the parameters.

1.1.2 Prediction

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

$$\hat{y} = \operatorname{argmax}_y \log p(\mathbf{x}, y; \boldsymbol{\mu}, \boldsymbol{\phi}) \quad [1.14]$$

$$= \operatorname{argmax}_y \log p(\mathbf{x} \mid y; \boldsymbol{\phi}) + \log p(y; \boldsymbol{\mu}) \quad [1.15]$$

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

$$\log p(\mathbf{x} \mid y; \boldsymbol{\phi}) + \log p(y; \boldsymbol{\mu}) = \log \left[B(\mathbf{x}) \prod_{j=1}^V \phi_{y,j}^{x_j} \right] + \log \mu_y \quad [1.16]$$

$$= \log B(\mathbf{x}) + \sum_j x_j \log \phi_{y,j} + \log \mu_y \quad [1.17]$$

$$= \log B(\mathbf{x}) + \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y), \quad [1.18]$$

where

$$\boldsymbol{\theta} = [\boldsymbol{\theta}^{(1)}; \boldsymbol{\theta}^{(2)}; \dots; \boldsymbol{\theta}^{(K)}] \quad [1.19]$$

$$\boldsymbol{\theta}^{(y)} = [\log \phi_{y,1}, \log \phi_{y,2}, \dots, \log \phi_{y,V}, \log \mu_y]^\top \quad [1.20]$$

The feature function $\mathbf{f}(\mathbf{x}, y)$ 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, and Figure 1.1). This construction ensures that the inner product $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y)$ only activates the features whose weights are in $\boldsymbol{\theta}^{(y)}$. These features and weights are all 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 $\langle \mathbf{x}, y \rangle$ into the computation of a vector inner product.

1.1.3 Estimation

The parameters of the categorical and multinomial distributions have a simple interpretation: they are vectors of expected frequencies for each possible event. Based on this interpretation, it is tempting to set the parameters empirically, as

$$\phi_{y,j} = \frac{\sum_{i:y^{(i)}=y} x_j^{(i)}}{\sum_{j'=1}^V \sum_{i:y^{(i)}=y} x_{j'}^{(i)}} = \frac{\text{count}(y, j)}{\sum_{j'=1}^V \text{count}(y, j')}, \quad [1.21]$$

where $\text{count}(y, j)$ refers to the count of word j in documents with label y .

Equation 1.21 defines the **relative frequency estimator** for ϕ . It can be justified more rigorously as a **maximum likelihood estimate**. Based on the generative model that we have defined, the log-likelihood of the entire corpus is:

$$\mathcal{L}(\phi, \mu) = \sum_{i=1}^N \log p_{\text{mult}}(\mathbf{x}^{(i)}; \phi_{y^{(i)}}) + \log p_{\text{cat}}(y^{(i)}; \mu), \quad [1.22]$$

which is now written as a function \mathcal{L} of the parameters ϕ and μ . Let's continue to focus on the parameters ϕ . Since $p(y)$ is constant with respect to ϕ , we can drop it:

$$\mathcal{L}(\phi) = \sum_{i=1}^N \log p_{\text{mult}}(\mathbf{x}^{(i)}; \phi_{y^{(i)}}) = \sum_{i=1}^N \log B(\mathbf{x}^{(i)}) + \sum_{j=1}^V x_j^{(i)} \log \phi_{y^{(i)},j}, \quad [1.23]$$

where $B(\mathbf{x}^{(i)})$ is constant with respect to ϕ .

We would now like to optimize the log-likelihood \mathcal{L} , by taking derivatives with respect to ϕ . But before we can do that, we have to deal with a set of constraints:

$$\sum_{j=1}^V \phi_{y,j} = 1 \quad \forall y \quad [1.24]$$

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

$$\ell(\phi_y) = \sum_{i:y^{(i)}=y} \sum_{j=1}^V x_j^{(i)} \log \phi_{y,j} - \lambda \left(\sum_{j=1}^V \phi_{y,j} - 1 \right). \quad [1.25]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

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

$$\frac{\partial \ell(\phi_y)}{\partial \phi_{y,j}} = \sum_{i:y^{(i)}=y} x_j^{(i)} / \phi_{y,j} - \lambda \quad [1.26]$$

$$\lambda \phi_{y,j} = \sum_{i:y^{(i)}=y} x_j^{(i)} \quad [1.27]$$

$$\phi_{y,j} \propto \sum_{i:y^{(i)}=y} x_j^{(i)} = \sum_{i=1}^N \delta(y^{(i)} = y) x_j^{(i)} = \text{count}(y, j), \quad [1.28]$$

where $\delta(y^{(i)} = y)$ is a function that returns one if $y^{(i)} = y$, and zero otherwise. Equation 1.28 shows three different notations for 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 $\sum_{j=1}^V \phi_{y,j} = 1$, which arises because ϕ_y represents a vector of probabilities for each word in the vocabulary. This constraint leads to an exact solution,

$$\phi_{y,j} = \frac{\text{count}(y, j)}{\sum_{j'=1}^V \text{count}(y, j')}. \quad [1.29]$$

This is equal to the relative frequency estimator from Equation 1.21. A similar derivation gives $\mu_y \propto \sum_{i=1}^N \delta(y^{(i)} = y)$.

1.1.4 Smoothing and MAP estimation

With text data, there are likely to be pairs of labels and words that never appear in the training set, leaving $\phi_{y,j} = 0$. For example, the word *Bayesian* may have never yet appeared in a spam email. But choosing a value of $\phi_{\text{SPAM}, \text{Bayesian}} = 0$ would allow this single feature to completely veto a label, since $p(\text{SPAM} | \mathbf{x}) = 0$ 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 a “pseudocount” of α to each count, and then normalizing.

$$\phi_{y,j} = \frac{\alpha + \text{count}(y, j)}{V\alpha + \sum_{j'=1}^V \text{count}(y, j')} \quad [1.30]$$

This is called **Laplace smoothing**.⁹ The pseudocount α is a **hyperparameter**, because it controls the form of the log-likelihood function, which in turn drives the estimation of ϕ .

⁹Laplace smoothing has a nice Bayesian justification, in which the generative model is extended to include ϕ as a random variable. 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 (Mohri et al., 2012). Techniques for reducing variance typically increase the bias, leading to a **bias-variance tradeoff**.

- Unbiased classifiers may **overfit** the training data, yielding poor performance on unseen data.
- But if the smoothing is too large, the resulting classifier can **underfit** instead. In the limit of $\alpha \rightarrow \infty$, there is zero variance: you get the same classifier, regardless of the data. However, the bias is likely to be large.

1.1.5 Setting hyperparameters

How should we choose the best value of hyperparameters like α ? Maximum likelihood will not work: the maximum likelihood estimate of α on the training set will always be $\alpha = 0$. In many cases, what we really want is **accuracy**: the number of correct predictions, divided by the total number of predictions. (Other measures of classification performance are discussed in § 3.4.) As we will see, it is hard to optimize for accuracy directly. But for scalar hyperparameters like α , we can employ a simple heuristic called **grid search**: try a set of values (e.g., $\alpha \in \{0.001, 0.01, 0.1, 1, 10\}$), compute the accuracy for each value, and choose the setting that maximizes the accuracy.

The goal is to tune α so that the classifier performs well on *unseen* data. For this reason, the data used for hyperparameter tuning should not overlap the training set, where very small values of α will be preferred. Instead, we hold out a **development set** (also called a **tuning set**) for hyperparameter selection. This development set may consist of a small fraction of the labeled data, such as 10%.

We will also want to predict the performance of our classifier on unseen data. To do this, we must hold out a separate subset of data, called the **test set**. It is critical that the test set not overlap with either the training or development sets, or else we risk overestimating the performance that the classifier will achieve on unlabeled data in the future. The test set should also not be used when making modeling decisions, such as the form of the feature function, the size of the vocabulary, and so on (these decisions are reviewed in chapter 3.) The ideal practice is to use the test set only once — otherwise, the test set is used to guide the classifier design, and test set accuracy will diverge from accuracy on truly unseen data. Because annotated data is expensive, this ideal can be hard to follow in practice, and many test sets have been used for decades. But in some high-impact applications like machine translation and information extraction, new test sets are released every year.

When only a small amount of labeled data is available, the test set accuracy can be unreliable. K -fold **cross-validation** is one way to cope with this scenario: the labeled

data is divided into K folds, and each fold acts as the test set, while training on the other folds. The test set accuracies are then aggregated. In the extreme, each fold is a single data point; this is called **leave-one-out** cross-validation. To perform hyperparameter tuning in the context of cross-validation, another fold can be used for grid search. It is important not to repeatedly evaluate the cross-validated accuracy while making design decisions about the classifier, or you will overstate the accuracy on truly unseen data.

1.2 Discriminative learning

Naïve Bayes is easy to work with: the weights can be estimated in closed form, and the probabilistic interpretation makes it relatively easy to extend. However, the assumption that features are independent can seriously limit its accuracy. Thus far, we have defined the **feature function** $f(x, y)$ so that it corresponds to bag-of-words features: one feature per word in the vocabulary. In natural language, bag-of-words features violate the assumption of conditional independence — for example, the probability that a document will contain the word *naïve* is surely higher given that it also contains the word *Bayes* — but this violation is relatively mild.

However, good performance on text classification often requires features that are richer than the bag-of-words:

- To better handle out-of-vocabulary terms, we want features that apply to multiple words, such as prefixes and suffixes (e.g., *anti-*, *un-*, *-ing*) and capitalization.
- We also want n -**gram** features that apply to multi-word units: **bigrams** (e.g., *not good*, *not bad*), **trigrams** (e.g., *not so bad*, *lacking any decency*, *never before imagined*), and beyond.

These features flagrantly violate the Naïve Bayes independence assumption. Consider what happens if we add a prefix feature. Under the Naïve Bayes assumption, we make the following approximation:¹⁰

$$\Pr(\text{word} = \textit{unfit}, \text{prefix} = \textit{un-} \mid y) \approx \Pr(\text{prefix} = \textit{un-} \mid y) \times \Pr(\text{word} = \textit{unfit} \mid y).$$

To test the quality of the approximation, we can manipulate the left-hand side by applying the chain rule,

$$\Pr(\text{word} = \textit{unfit}, \text{prefix} = \textit{un-} \mid y) = \Pr(\text{prefix} = \textit{un-} \mid \text{word} = \textit{unfit}, y) \quad [1.31]$$

$$\times \Pr(\text{word} = \textit{unfit} \mid y) \quad [1.32]$$

¹⁰The notation $\Pr(\cdot)$ refers to the probability of an event, and $p(\cdot)$ refers to the probability density or mass for a random variable (see Appendix A).

But $\Pr(\text{prefix} = \text{un-} \mid \text{word} = \text{unfit}, y) = 1$, since *un-* is guaranteed to be the prefix for the word *unfit*. Therefore,

$$\begin{aligned} \Pr(\text{word} = \text{unfit}, \text{prefix} = \text{un-} \mid y) &= 1 \times \Pr(\text{word} = \text{unfit} \mid y) \quad [1.33] \\ &\gg \Pr(\text{prefix} = \text{un-} \mid y) \times \Pr(\text{word} = \text{unfit} \mid y), \quad [1.34] \end{aligned}$$

because the probability of any given word starting with the prefix *un-* is much less than one. Naïve Bayes will systematically underestimate the true probabilities of conjunctions of positively correlated features. To use such features, we need learning algorithms that do not rely on an independence assumption.

The origin of the Naïve Bayes independence assumption is the learning objective, $p(\mathbf{x}^{(1:N)}, y^{(1:N)})$, which requires modeling the probability of the observed text. In classification problems, we are always given \mathbf{x} , and are only interested in predicting the label y , so it seems unnecessary to model the probability of \mathbf{x} . **Discriminative learning** algorithms focus on the problem of predicting y , and do not attempt to model the probability of the text \mathbf{x} .

1.2.1 Perceptron

In Naïve Bayes, 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 3, is one way to do this.

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 is an **online learning** algorithm, since the classifier weights change after every example. This is different from Naïve Bayes, which computes corpus statistics and then sets the weights in a single operation — Naïve Bayes is a **batch learning** algorithm. Algorithm 3 is vague about when this online learning procedure terminates. We will return to this issue shortly.

The perceptron algorithm may seem like a cheap heuristic: Naïve Bayes has a solid foundation in probability, but now we are just adding and subtracting constants from the weights every time there is a mistake. Will this really work? In fact, there is some nice theory for the perceptron, based on the concept of **linear separability**:

Definition 1 (Linear separability). *The dataset $\mathcal{D} = \{\mathbf{x}^{(1:N)}, y^{(1:N)}\}$ is linearly separable iff there exists some weight vector $\boldsymbol{\theta}$ and some **margin** ρ such that for every instance $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$, the inner product of $\boldsymbol{\theta}$ and the feature function for the true label, $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$, is at least ρ greater than inner product of $\boldsymbol{\theta}$ and the feature function for every other possible label, $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')$.*

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 3 Perceptron learning algorithm

```

1: procedure PERCEPTRON( $\mathbf{x}^{(1:N)}, y^{(1:N)}$ )
2:    $t \leftarrow 0$ 
3:    $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ 
4:   repeat
5:      $t \leftarrow t + 1$ 
6:     Select an instance  $i$ 
7:      $\hat{y} \leftarrow \operatorname{argmax}_y \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$ 
8:     if  $\hat{y} \neq y^{(i)}$  then
9:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$ 
10:    else
11:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
12:  until tired
13:  return  $\boldsymbol{\theta}^{(t)}$ 

```

$$\exists \boldsymbol{\theta}, \rho > 0 : \forall \langle \mathbf{x}^{(i)}, y^{(i)} \rangle \in \mathcal{D}, \quad \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) \geq \rho + \max_{y' \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y'). \quad [1.35]$$

Linear separability is important because of the following guarantee: if your data is linearly separable, then the perceptron algorithm will find a separator (Novikoff, 1962).¹¹ So while the perceptron may seem heuristic, it is guaranteed to succeed, if the learning problem is easy enough.

How useful is this proof? Minsky and Papert (1969) famously proved that the simple logical function of *exclusive-or* is not separable, and that a perceptron is therefore incapable of learning this function. But this is not just an issue for the perceptron: any linear classification algorithm, including Naïve Bayes, will fail on this task. In natural language classification problems usually involve high dimensional feature spaces, with thousands or millions of features. For these problems, it is very likely that the training data is indeed separable. And even if the data is not separable, it is still possible to place an upper bound on the number of errors that the perceptron algorithm will make (Freund and Schapire, 1999).

1.2.2 Averaged perceptron

The perceptron iterates over the data repeatedly — until “tired”, as described in Algorithm 3. If the data is linearly separable, the perceptron will eventually find a separator,

¹¹It is also possible to prove an upper bound on the number of training iterations required to find the separator. Proofs like this are part of the field of **statistical learning theory** (Mohri et al., 2012).

Algorithm 4 Averaged perceptron learning algorithm

```

1: procedure AVG-PERCEPTRON( $\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)}$ )
2:    $t \leftarrow 0$ 
3:    $\boldsymbol{\theta}^{(0)} \leftarrow \mathbf{0}$ 
4:   repeat
5:      $t \leftarrow t + 1$ 
6:     Select an instance  $i$ 
7:      $\hat{y} \leftarrow \operatorname{argmax}_y \boldsymbol{\theta}^{(t-1)} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$ 
8:     if  $\hat{y} \neq y^{(i)}$  then
9:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} + \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$ 
10:    else
11:       $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)}$ 
12:       $\mathbf{m} \leftarrow \mathbf{m} + \boldsymbol{\theta}^{(t)}$ 
13:    until tired
14:     $\bar{\boldsymbol{\theta}} \leftarrow \frac{1}{t} \mathbf{m}$ 
15:  return  $\bar{\boldsymbol{\theta}}$ 

```

and we can stop once all training instances are classified correctly. But if the data is not linearly separable, the perceptron can *thrash* between two or more weight settings, never converging. In this case, how do we know that we can stop training, and how should we choose the final weights? An effective practical solution is to *average* the perceptron weights across all iterations.

This procedure is shown in Algorithm 4. The learning algorithm is nearly identical, but we also maintain a vector of the sum of the weights, \mathbf{m} . At the end of the learning procedure, we divide this sum by the total number of updates t , to compute the average weights, $\bar{\boldsymbol{\theta}}$. These average weights are then used for prediction. In the algorithm sketch, the average is computed from a running sum, $\mathbf{m} \leftarrow \mathbf{m} + \boldsymbol{\theta}$. However, this is inefficient, because it requires $|\boldsymbol{\theta}|$ operations to update the running sum. When $\mathbf{f}(\mathbf{x}, y)$ is sparse, $|\boldsymbol{\theta}| \gg |\mathbf{f}(\mathbf{x}, y)|$ for any individual (\mathbf{x}, y) . This means that computing the running sum will be much more expensive than computing of the update to $\boldsymbol{\theta}$ itself, which requires only $2 \times |\mathbf{f}(\mathbf{x}, y)|$ operations. One of the exercises is to sketch a more efficient algorithm for computing the averaged weights.

Even if the data is not separable, the averaged weights will eventually converge. One possible stopping criterion is to check the difference between the average weight vectors after each pass through the data: if the norm of the difference falls below some predefined threshold, we can stop iterating. Another stopping criterion is to hold out some data, and to measure the predictive accuracy on this heldout data. When the accuracy on the heldout data starts to decrease, the learning algorithm has begun to **overfit** the training set. At this point, it is probably best to stop; this stopping criterion is known as **early**

stopping.

Generalization is the ability to make good predictions on instances that are not in the training data. Averaging can be proven to improve generalization, by computing an upper bound on the generalization error (Freund and Schapire, 1999; Collins, 2002).

1.3 Loss functions and large-margin classification

Naïve Bayes chooses the weights θ by maximizing the joint log-likelihood $\log p(\mathbf{x}^{(1:N)}, y^{(1:N)})$. Optimization problems are generally formulated as minimization of a **loss function**. We can trivially reformulate maximum likelihood in this way, by defining the loss function to be the *negative* log-likelihood:

$$\log p(\mathbf{x}^{(1:N)}, y^{(1:N)}; \theta) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}, y^{(i)}; \theta) \quad [1.36]$$

$$\ell_{\text{NB}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = -\log p(\mathbf{x}^{(i)}, y^{(i)}; \theta) \quad [1.37]$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmin}} \sum_{i=1}^N \ell_{\text{NB}}(\theta; \mathbf{x}^{(i)}, y^{(i)}), \quad [1.38]$$

This minimization problem is identical to maximum-likelihood estimation. Loss functions provide a general framework for comparing machine learning objectives. For example, an alternative loss function is the **zero-one loss**,

$$\ell_{0-1}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \begin{cases} 0, & y^{(i)} = \operatorname{argmax}_y \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) \\ 1, & \text{otherwise} \end{cases} \quad [1.39]$$

This loss function is proportional to the error rate of the classifier on the training data. Since a low error rate is often the ultimate goal of classification, this may seem ideal. But the zero-one loss is **non-convex**,¹² which means that it is combinatorially difficult to optimize.

The perceptron optimizes the following loss function:

$$\ell_{\text{PERCEPTRON}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \max_{y \in \mathcal{Y}} \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) - \theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}), \quad [1.40]$$

which is a **hinge loss** with the hinge point at zero. When $\hat{y} = y^{(i)}$, the loss is zero; otherwise, it increases linearly with the gap between the score for the predicted label \hat{y} and the score for the true label $y^{(i)}$.

¹²A 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 many effective techniques for optimizing convex functions (Boyd and Vandenberghe, 2004).

To see why this is the loss function optimized by the perceptron, take the derivative with respect to θ ,

$$\frac{\partial}{\partial \theta} \ell_{\text{PERCEPTRON}}(\theta; \mathbf{x}^{(i)}, y^{(i)}) = \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}) - \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}). \quad [1.41]$$

At each instance perceptron algorithm takes a step of magnitude one in the opposite direction of this **gradient**, $\nabla_{\theta} \ell_{\text{PERCEPTRON}} = \frac{\partial}{\partial \theta} \ell_{\text{PERCEPTRON}}(\theta; \mathbf{x}^{(i)}, y^{(i)})$. As we will see in § 1.5, this is an example of the optimization algorithm **stochastic gradient descent**, applied to the objective in Equation 1.40.

Breaking ties with subgradient descent Careful readers will notice the tacit assumption that there is a unique \hat{y} that maximizes $\theta \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$. What if there are two or more labels that maximize this function? Consider binary classification: if the maximizer is $y^{(i)}$, then the gradient is zero, and so is the perceptron update; if the maximizer is $\hat{y} \neq y^{(i)}$, then the update is the difference $\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$. The underlying issue is that the perceptron loss is not **smooth**, because the first derivative has a discontinuity at the hinge point, where the score for the true label $y^{(i)}$ is equal to the score for some other label \hat{y} . At this point, there is no unique gradient; rather, there is a set of **subgradients**. A vector \mathbf{v} is a subgradient of the function g at \mathbf{u}_0 iff $g(\mathbf{u}) - g(\mathbf{u}_0) \geq \mathbf{v} \cdot (\mathbf{u} - \mathbf{u}_0)$ for all \mathbf{u} . Graphically, this defines the set of hyperplanes that include $g(\mathbf{u}_0)$ and do not intersect g at any other point. As we approach the hinge point from the left, the gradient is $\mathbf{f}(\mathbf{x}, \hat{y}) - \mathbf{f}(\mathbf{x}, y)$; as we approach from the right, the gradient is $\mathbf{0}$. At the hinge point, the subgradients include all vectors that are bounded by these two extremes. In subgradient descent, *any* subgradient can be used (Bertsekas, 2012). Since both $\mathbf{0}$ and $\mathbf{f}(\mathbf{x}, \hat{y}) - \mathbf{f}(\mathbf{x}, y)$ are subgradients at the hinge point, either one can be used in the perceptron update.

Perceptron versus Naïve Bayes The perceptron loss function has some pros and cons with respect to the negative log-likelihood loss implied by Naïve Bayes.

- Both ℓ_{NB} and $\ell_{\text{PERCEPTRON}}$ are convex, making them relatively easy to optimize. However, ℓ_{NB} can be optimized in closed form, while $\ell_{\text{PERCEPTRON}}$ requires iterating over the dataset multiple times.
- ℓ_{NB} can suffer **infinite** loss on a single example, since the logarithm of zero probability is negative infinity. Naïve Bayes will therefore overemphasize some examples, and underemphasize others.
- $\ell_{\text{PERCEPTRON}}$ treats all correct answers equally. Even if θ only gives the correct answer by a tiny margin, the loss is still zero.

(c) Jacob Eisenstein 2014-2017. Work in progress.

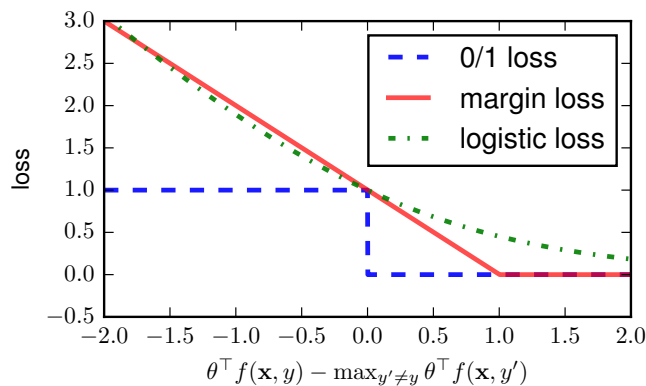


Figure 1.2: Margin, zero-one, and logistic loss functions.[**todo: redo x-axis labels, add hinge loss with hinge point at zero.**]

1.3.1 Large margin classification

This last comment suggests a potential problem with the perceptron. Suppose a test example is very close to a training example, but not identical. If the classifier only gets the correct answer on the training example by a small margin, then it may get the test instance wrong. To formalize this intuition, define the **margin** as,

$$\gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \max_{y \neq y^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y). \quad [1.42]$$

The margin represents the difference between the score for the correct label $y^{(i)}$, and the score for the highest-scoring label. The intuition behind **large margin classification** is that it is not enough just to label the training data correctly — the correct label should be separated from other labels by a comfortable margin. This idea can be encoded into a loss function,

$$\ell_{\text{MARGIN}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = \begin{cases} 0, & \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) \geq 1, \\ 1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}), & \text{otherwise} \end{cases} \quad [1.43]$$

$$= \left(1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)})\right)_+, \quad [1.44]$$

where $(x)_+ = \max(0, x)$. The loss is zero if there is a margin of at least 1 between the score for the true label and the best-scoring alternative \hat{y} . This is almost identical to the perceptron loss, but the hinge point is shifted to the right, as shown in Figure 1.2. The margin loss is a convex upper bound on the zero-one loss.

(c) Jacob Eisenstein 2014-2017. Work in progress.

1.3.2 Support vector machines

If a dataset is linearly separable, then there is some hyperplane θ that correctly classifies all training instances with margin ρ (by Definition 1). This margin can be increased to any desired value by multiplying the weights by a constant. Now, for any datapoint $\langle \mathbf{x}^{(i)}, y^{(i)} \rangle$, the geometric distance to the separating hyperplane is given by $\frac{\gamma(\theta; \mathbf{x}^{(i)}, y^{(i)})}{\|\theta\|_2}$, where the denominator is the norm of the weights, $\|\theta\|_2 = \sqrt{\sum_j \theta_j^2}$. The geometric distance is sometimes called the **geometric margin**, in contrast to the **functional margin** $\gamma(\theta; \mathbf{x}^{(i)}, y^{(i)})$. Both are shown in Figure 1.3. The geometric margin is a good measure of the robustness of the separator: if the functional margin is large, but the norm $\|\theta\|_2$ is also large, then a small change in $\mathbf{x}^{(i)}$ could cause it to be misclassified. We therefore seek to maximize the minimum geometric margin, subject to the constraint that the functional margin is at least one:

$$\begin{aligned} \max_{\theta} . \quad & \min_i . \quad \frac{\gamma(\theta; \mathbf{x}^{(i)}, y^{(i)})}{\|\theta\|_2} \\ \text{s.t.} \quad & \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}) \geq 1, \quad \forall i. \end{aligned} \tag{1.45}$$

This is a **constrained optimization** problem, where the second line describes constraints on the space of possible solutions θ . In this case, the constraint is that the functional margin always be at least one, and the objective is that the minimum geometric margin be as large as possible.

Any scaling factor on θ will cancel in the numerator and denominator of the geometric margin. This means that if the data is linearly separable at ρ , we can increase this margin to 1 by rescaling θ . We therefore need only minimize the denominator $\|\theta\|_2$, subject to the constraint on the functional margin. The minimizer of $\|\theta\|_2$ is also the minimizer of $\frac{1}{2}\|\theta\|_2^2 = \frac{1}{2} \sum_j \theta_j^2$, which is easier to work with. This gives the optimization problem,

$$\begin{aligned} \min_{\theta} . \quad & \frac{1}{2} \|\theta\|_2^2 \\ \text{s.t.} \quad & \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}) \geq 1, \quad \forall i. \end{aligned} \tag{1.46}$$

This optimization problem is a **quadratic program**: the objective is a quadratic function of the parameters, and the constraints are all linear inequalities. The resulting classifier is better known as the **support vector machine**. The name derives from one of the solutions, which is to incorporate the constraints through Lagrange variables $\alpha_i \geq 0, i \in \{1, 2, \dots, N\}$. The instances for which $\alpha_i > 0$ are the **support vectors**; other instances are irrelevant to the classification boundary.

(c) Jacob Eisenstein 2014-2017. Work in progress.

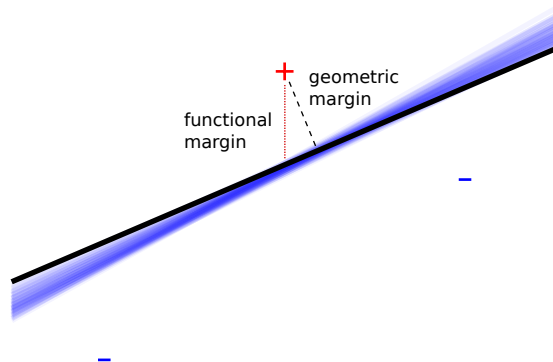


Figure 1.3: Functional and geometric margins for a binary classification problem. All separators that satisfy the margin constraint are shown. The separator with the largest geometric margin is shown in bold.

1.3.3 Slack variables

If a dataset is not linearly separable, then there is no θ that satisfies the margin constraint. To add more flexibility, we introduce a set of **slack variables** $\xi_i \geq 0$. Instead of requiring that the functional margin be greater than or equal to one, we require that it be greater than or equal to $1 - \xi_i$. Ideally there would not be any slack, so the slack variables are penalized in the objective function:

$$\begin{aligned} \min_{\theta, \xi} \quad & \frac{1}{2} \|\theta\|_2^2 + C \sum_{i=1}^N \xi_i \\ \text{s.t.} \quad & \gamma(\theta; \mathbf{x}^{(i)}, y^{(i)}) + \xi_i \geq 1, \quad \forall i \\ & \xi_i \geq 0, \quad \forall i. \end{aligned} \tag{1.47}$$

The hyperparameter C controls the tradeoff between violations of the margin constraint and the preference for a low norm of θ . As $C \rightarrow \infty$, slack is infinitely expensive, and there is only a solution if the data is separable. As $C \rightarrow 0$, slack becomes free, and there is a trivial solution at $\theta = \mathbf{0}$. Thus, C plays a similar role to the smoothing parameter in Naïve Bayes (§ 1.1.4), trading off between a close fit to the training data and better generalization. Like the smoothing parameter of Naïve Bayes, C must be set by the user, typically by maximizing performance on a heldout development set.

To solve the constrained optimization problem defined in Equation 1.47, we can first

solve for the slack variables,

$$\xi_i \geq (1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}))_+. \quad [1.48]$$

The inequality is tight, because the slack variables are penalized in the objective, and there is no advantage to increasing them beyond the minimum value (Ratliff et al., 2007; Smith, 2011). The problem can therefore be transformed into the unconstrained optimization,

$$\min_{\boldsymbol{\theta}} \quad \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2 + \sum_{i=1}^N (1 - \gamma(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}))_+, \quad [1.49]$$

where each ξ_i has been substituted by the right-hand side of Equation 1.48, and the factor of C on the slack variables has been replaced by an equivalent factor of $\lambda = \frac{1}{C}$ on the norm of the weights.

Now define the **cost** of a classification error as,¹³

$$c(y^{(i)}, \hat{y}) = \begin{cases} 1, & y^{(i)} \neq \hat{y} \\ 0, & \text{otherwise.} \end{cases} \quad [1.50]$$

We can then rewrite Equation 1.49,

$$\min_{\boldsymbol{\theta}} \quad \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2 + \sum_{i=1}^N \left(\max_{y \in \mathcal{Y}} \left(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) + c(y^{(i)}, y) \right) - \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) \right)_+, \quad [1.51]$$

maximizing over all $y \in \mathcal{Y}$, in search of labels that are both *high scoring*, as measured by $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)$, and *wrong*, as measured by $c(y^{(i)}, y)$. This maximization is known as **cost-augmented inference**, because it augments the maximization objective to favor high-cost predictions. If the highest-scoring label is $y = y^{(i)}$, then the margin constraint is satisfied, and the loss for this instance is zero. Note that cost-augmentation is not applied when making predictions on unseen data.

Differentiating Equation 1.51 with respect to the weights gives,

$$\nabla_{\boldsymbol{\theta}} L_{\text{SVM}} = \lambda \boldsymbol{\theta} + \sum_{i=1}^N \mathbf{f}(\mathbf{x}^{(i)}, \hat{y}) - \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) \quad [1.52]$$

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) + c(y^{(i)}, y), \quad [1.53]$$

where L_{SVM} refers to minimization objective in Equation 1.51. This gradient is very similar to the perceptron update. One difference is the additional term $\lambda \boldsymbol{\theta}$, which **regularizes** the

¹³We can also define specialized cost functions that heavily penalize especially undesirable errors (Tsochantaridis et al., 2004). This issue is revisited in chapter 6.

weights towards $\mathbf{0}$. The other difference is the cost $c(y^{(i)}, y)$, which is added to $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y)$ when choosing \hat{y} during training. This term derives from the margin constraint: large margin classifiers learn not only from instances that are incorrectly classified, but also from instances for which the correct classification decision was not sufficiently confident.

1.4 Logistic regression

Thus far, we have seen two broad classes of learning algorithms. Naïve Bayes is a probabilistic method, where learning is equivalent to estimating a joint probability distribution. The perceptron and support vector machine are discriminative, error-driven algorithms: the learning objective is closely related to the number of errors on the training data. Probabilistic and error-driven approaches each have advantages: probability makes it possible to quantify uncertainty about the predicted labels, but the probability model of Naïve Bayes brings unrealistic independence assumptions that limit the features that can be used.

Logistic regression combines advantages of discriminative and probabilistic classifiers. Unlike Naïve Bayes, which starts from the **joint probability** $p(\mathbf{x}, y)$, logistic regression defines the desired **conditional probability** $p(y | \mathbf{x})$ directly. Think of $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y)$ as a scoring function for the compatibility of the base features \mathbf{x} and the label y . To convert this score into a probability, we first exponentiate, obtaining $\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y))$, which is guaranteed to be non-negative. Next, we normalize, dividing over all possible labels $y' \in \mathcal{Y}$. The resulting conditional probability is defined as,

$$p(y | \mathbf{x}; \boldsymbol{\theta}) = \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y))}{\sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y'))}. \quad [1.54]$$

Given a dataset $\mathcal{D} = \{\mathbf{x}^{(1:N)}, y^{(1:N)}\}$, the weights $\boldsymbol{\theta}$ are estimated by **maximum conditional likelihood**,

$$\log p(\mathbf{y}^{(1:N)} | \mathbf{x}^{(1:N)}; \boldsymbol{\theta}) = \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \quad [1.55]$$

$$= \sum_{i=1}^N \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \log \sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')). \quad [1.56]$$

The final line is obtained by plugging in Equation 1.54 and taking the logarithm.¹⁴ Inside

¹⁴The log-sum-exp term is a common pattern in machine learning. It is numerically unstable, because it will underflow if the inner product is small, and overflow if the inner product is large. Scientific computing libraries usually 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.

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

$$\ell_{\text{LOGREG}}(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) = -(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \log \sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y'))) \quad [1.57]$$

The logistic loss is shown in Figure 1.2. A key difference from the zero-one and hinge losses is that logistic loss is never exactly zero. This means that the objective function can always be improved by choosing the correct label with more confidence.

1.4.1 Regularization

As with the margin-based algorithms described in § 1.3, we can obtain better generalization by penalizing the norm of $\boldsymbol{\theta}$. This is done by adding a term of $\frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2$ to the minimization objective. This is called L_2 regularization, because $\|\boldsymbol{\theta}\|_2^2$ is the squared L_2 norm of the vector $\boldsymbol{\theta}$. L_2 regularization can be viewed as placing a zero-mean Gaussian prior distribution on each term of $\boldsymbol{\theta}$, because the log-likelihood under a zero-mean Gaussian is,

$$\log N(\theta_j; 0, \sigma^2) \propto -\frac{1}{2\sigma^2} \theta_j^2, \quad [1.58]$$

so that $\lambda = \frac{1}{\sigma^2}$.

As in the support vector machine classifier, L_2 regularization forces the estimator to trade performance on the training data for a smaller norm of the weights, and this can help to prevent overfitting. Regularization is generally considered to be essential to estimating high-dimensional models. 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 $\boldsymbol{\theta}_{(j, y^{(i)})} \rightarrow \infty$ and $\boldsymbol{\theta}_{(j, \tilde{y} \neq y^{(i)})} \rightarrow -\infty$, where (j, y) is the index of feature associated with $x_j^{(i)}$ and label y in $\mathbf{f}(\mathbf{x}^{(i)}, y)$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

1.4.2 Gradients

For a single example, here is the logistic loss and its gradient:

$$\ell_{\text{LOGREG}} = -\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \log \sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')) \quad [1.59]$$

$$\frac{\partial \ell}{\partial \boldsymbol{\theta}} = -\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \frac{1}{\sum_{y'' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y''))} \times \sum_{y' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y')) \times \mathbf{f}(\mathbf{x}^{(i)}, y') \quad [1.60]$$

$$= -\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \sum_{y' \in \mathcal{Y}} \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y'))}{\sum_{y'' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y''))} \times \mathbf{f}(\mathbf{x}^{(i)}, y') \quad [1.61]$$

$$= -\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + \sum_{y' \in \mathcal{Y}} p(y' | \mathbf{x}^{(i)}; \boldsymbol{\theta}) \times \mathbf{f}(\mathbf{x}^{(i)}, y') \quad [1.62]$$

$$= -\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) + E_{y|\mathbf{x}^{(i)}}[\mathbf{f}(\mathbf{x}^{(i)}, y)]. \quad [1.63]$$

The final step employs the definition of a conditional expectation (§ A.5). The gradient of the logistic loss is equal to the difference between the observed feature counts $\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})$ and the expected counts under the current model, $E_{y|\mathbf{x}^{(i)}}[\mathbf{f}(\mathbf{x}^{(i)}, y)]$. When these two vectors are equal for a single instance, there is nothing more to learn from it; when they are equal in sum over the entire dataset, there is nothing more to learn from the dataset as a whole. The gradient of the hinge loss is nearly identical, but it involves the features of the predicted label under the current model, $\mathbf{f}(\mathbf{x}^{(i)}, \hat{y})$, rather than the expected features $E_{y|\mathbf{x}^{(i)}}[\mathbf{f}(\mathbf{x}^{(i)}, y)]$ under the conditional distribution $p(y | \mathbf{x}; \boldsymbol{\theta})$.

The regularizer contributes $\lambda \boldsymbol{\theta}$ to the overall gradient:

$$L_{\text{LOGREG}} = \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2 - \sum_{i=1}^N \left(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \log \sum_{y' \in \mathcal{Y}} \exp \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y') \right) \quad [1.64]$$

$$\nabla_{\boldsymbol{\theta}} L_{\text{LOGREG}} = \lambda \boldsymbol{\theta} - \sum_{i=1}^N \left(\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - E_{y|\mathbf{x}}[\mathbf{f}(\mathbf{x}^{(i)}, y)] \right). \quad [1.65]$$

1.5 Optimization

Each of the classification algorithms in this chapter can be viewed as an optimization problem:

- In Naïve Bayes, the objective is the joint likelihood $\log p(\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)})$. Maximum likelihood estimation yields a closed-form solution for $\boldsymbol{\theta}$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- In the support vector machine, the objective is,

$$L_{\text{SVM}} = \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2 + \sum_{i=1}^N (\max_{y \in \mathcal{Y}} (\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y) + c(y^{(i)}, y)) - \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}))_+, \quad [1.66]$$

There is no closed-form solution, but the objective is convex. The perceptron algorithm minimizes a similar objective.

- In logistic regression, the objective is the regularized negative log-likelihood,

$$L_{\text{LOGREG}} = \frac{\lambda}{2} \|\boldsymbol{\theta}\|_2^2 - \sum_{i=1}^N \left(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y^{(i)}) - \log \sum_{y \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}^{(i)}, y)) \right) \quad [1.67]$$

Again, there is no closed-form solution, but the objective is convex.

These learning algorithms are distinguished by *what* is being optimized, rather than *how* the optimal weights are found. This decomposition is an essential feature of contemporary machine learning. The domain expert's job is to design an objective function — or more generally, a **model** of the problem. If the model has certain characteristics, then generic optimization algorithms can be used to find the solution. In particular, if an objective function is *differentiable*, then gradient-based optimization can be employed; if it is also *convex*, then gradient-based optimization is guaranteed to find the globally optimal solution. The support vector machine and logistic regression have both of these properties, and so are amenable to generic **convex optimization** techniques (Boyd and Vandenberghe, 2004).

1.5.1 Batch optimization

In **batch optimization**, each update to the weights is based on a computation involving the entire dataset. One such algorithm is **gradient descent**, which iteratively updates the weights,

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} - \eta^{(t)} \nabla_{\boldsymbol{\theta}} L \quad [1.68]$$

where $\nabla_{\boldsymbol{\theta}} L$ is the gradient computed over the entire training set, and $\eta^{(t)}$ is the **step size** at iteration t . If the objective L is a convex function of $\boldsymbol{\theta}$, then this procedure is guaranteed to terminate at the global optimum, as long as $\eta^{(t)}$ is small enough.

In practice, gradient descent can be slow to converge, as the gradient can become infinitesimally small. Faster convergence can be obtained by “second-order” Newton optimization, which incorporates the inverse of the **Hessian matrix**,

$$H_{i,j} = \frac{\partial^2 \ell}{\partial \theta_i \partial \theta_j} \quad [1.69]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

The size of the Hessian matrix is quadratic in the number of features. In the bag-of-words representation, this is usually too big to store, let alone invert. **Quasi-Network optimization** techniques maintain a low-rank approximation to the inverse of the Hessian matrix. Such techniques usually converge more quickly than gradient descent, while remaining computationally tractable even for large feature sets. A popular quasi-Newton algorithm is **L-BFGS** (Liu and Nocedal, 1989), which is implemented in many scientific computing environments, such as `scipy` and `Matlab`.

For any gradient-based technique, the user must set the learning rates $\eta^{(t)}$. Convergence proofs often use the schedule $\eta^{(t)} = \eta^{(0)} t^{-\alpha}$ for $\alpha \in [1, 2]$. In practice, $\eta^{(t)}$ may also be fixed to a small constant, like 10^{-3} . In either case, it is typical to try a few different values, and see which minimizes the objective most quickly.

1.5.2 Online optimization

Batch optimization computes the objective on the entire training set before making an update. This may be inefficient, because in many cases, a small number of training examples could point the learner in roughly the correct direction. **Online learning** algorithms make updates to the weights while iterating through the training data. The theoretical basis for this approach is a stochastic approximation to the true objective function,

$$\sum_{i=1}^N \ell(\theta; \mathbf{x}^{(i)}, y^{(i)}) \approx N \times \ell(\theta; \mathbf{x}^{(j)}, y^{(j)}), \quad \langle \mathbf{x}^{(j)}, y^{(j)} \rangle \sim \{\mathbf{x}^{(1:N)}, y^{(1:N)}\}, \quad [1.70]$$

where the instance $\langle \mathbf{x}^{(j)}, y^{(j)} \rangle$ is sampled at random from the full dataset.

In **stochastic gradient descent**, we randomly sample a single instances, and make an immediate update based on the local gradient. Note the similarity to the perceptron algorithm, which also updates the weights one instance at a time. In **minibatch** stochastic gradient descent, the gradient is computed over a small set of instances. A typical approach is to set the minibatch size so that the entire batch fits in memory on the GPU (Neubig et al., 2017b). It is then possible to speed up learning by parallelizing the computation of the gradient over each instance in the minibatch.

Algorithm 5 offers a generalized view of gradient descent. If the batcher returns a single batch with all instances, then the algorithm is standard gradient descent; if it returns N batches with one instance each, then it is stochastic gradient descent; otherwise, the batcher can return B minibatches, $1 < B < N$.

There are a number of more advanced techniques for online learning (Bottou et al., 2016). Many algorithms use an adaptive step size, which can be different for every feature (Duchi et al., 2011). Features that occur frequently are likely to be updated frequently, so it is best to use a small step size; rare features will be updated infrequently, so it is better to take larger steps. The **AdaGrad** (adaptive gradient) algorithm achieves this behavior by

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 5 Generalized gradient descent. The function `Batcher` partitions the training set into B batches such that each instance appears in exactly one batch. In gradient descent, $B = 1$; in stochastic gradient descent, $B = N$; in minibatch stochastic gradient descent, $1 < B < N$.

```

1: procedure GRADIENT-DESCENT( $\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)}, L, \eta^{(1:\infty)}, \text{Batcher}, T_{\max}$ )
2:    $\boldsymbol{\theta} \leftarrow \mathbf{0}$ 
3:    $t \leftarrow 0$ 
4:   repeat
5:      $\langle \mathbf{b}^{(1)}, \mathbf{b}^{(2)}, \dots, \mathbf{b}^{(B)} \rangle \leftarrow \text{Batcher}(N)$ 
6:     for  $n \in \{1, 2, \dots, B\}$  do
7:        $t \leftarrow t + 1$ 
8:        $\boldsymbol{\theta}^{(t)} \leftarrow \boldsymbol{\theta}^{(t-1)} - \eta^{(t)} \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^{(t-1)}; \mathbf{x}^{(b_1^{(n)}, b_2^{(n)}, \dots)}, \mathbf{y}^{(b_1^{(n)}, b_2^{(n)}, \dots)})$ 
9:       if Converged( $\boldsymbol{\theta}^{(1,2,\dots,t)}$ ) then
10:        return  $\boldsymbol{\theta}^{(t)}$ 
11:   until  $t \geq T_{\max}$ 
12:   return  $\boldsymbol{\theta}^{(t)}$ 

```

storing the sum of the squares of the gradients for each feature, and rescaling the learning rate by its inverse:

$$\mathbf{g}_t = \nabla_{\boldsymbol{\theta}} L(\boldsymbol{\theta}^{(t)}; \mathbf{x}^{(i)}, y^{(i)}) \quad [1.71]$$

$$\theta_j^{(t+1)} \leftarrow \theta_j^{(t)} - \frac{\eta^{(t)}}{\sqrt{\sum_{t'=1}^t g_{t,j}^2}} g_{t,j}, \quad [1.72]$$

where j iterates over features in $\mathbf{f}(\mathbf{x}, y)$.

In most cases, the number of active features for any instance is much smaller than the number of weights, $|\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})| \ll |\boldsymbol{\theta}|$. If so, the computation cost of online optimization will be dominated by the update from the regularization term, $\lambda \boldsymbol{\theta}$. The solution is to be “lazy”, updating each θ_j only as it is used. To implement lazy updating, store an additional parameter τ_j , which is the iteration at which θ_j was last updated. If θ_j is needed at time t , the $t - \tau_j$ regularization updates can be performed all at once. This strategy is described by Kummerfeld et al. (2015).

1.6 *Additional topics in classification

1.6.1 Sparsity-inducing regularization

In logistic regression and large-margin classification, generalization can be improved by regularizing the weights towards 0. The most popular regularizer is the squared L_2 norm

(c) Jacob Eisenstein 2014-2017. Work in progress.

$\|\theta\|_2^2$. However, we might prefer some other norm, such as $L_0 = \|\theta\|_0 = \sum_j \delta(\theta_j \neq 0)$, which applies a constant penalty for each non-zero weight. This norm 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 irrelevant features will be pruned away. Unfortunately, the L_0 norm is non-convex and non-differentiable; optimization under L_0 regularization is **NP-hard**, meaning that it can be solved efficiently only if $P=NP$ (Ge et al., 2011).

A useful alternative is the L_1 norm, which is equal to the sum of the absolute values of the weights, $\|\theta\|_1 = \sum_j |\theta_j|$. The L_1 norm is convex, and can be used as an approximation to L_0 (Tibshirani, 1996). Conveniently, 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**. The L_1 norm does not have a gradient at $\theta_j = 0$, so we must instead optimize the L_1 -regularized objective using **subgradient** methods. The associated stochastic sub-gradient descent algorithms are only somewhat more complex than conventional SGD; Sra et al. (2012) survey approaches for estimation under L_1 and other regularizers.

Gao et al. (2007) compare L_1 and L_2 regularization on a suite of NLP problems, finding that L_1 regularization generally gives similar accuracy to L_2 regularization, but that L_1 regularization produces models that are between ten and fifty times smaller, because more than 90% of the feature weights are set to zero.

1.6.2 Other views of logistic regression

In binary classification, we can dispense with the feature function, and choose y based on the inner product of $\theta \cdot x$. We can obtain a conditional probability $p(y \mid x; \theta)$ by passing this inner product through a **logistic function**,

$$\sigma(a) \triangleq \frac{\exp(a)}{1 + \exp(a)} = (1 + \exp(-a))^{-1} \quad [1.73]$$

$$p(y \mid x; \theta) = \sigma(\theta \cdot x). \quad [1.74]$$

This is the origin of the name **logistic regression**. Logistic regression can be viewed as part of a larger family of **generalized linear models** (GLMs), in which various other **link functions** convert between the inner product $\theta \cdot x$ and the parameter of a probability distribution $p(y \mid x)$.

In the early NLP literature, logistic regression is frequently called **maximum entropy** classification (Berger et al., 1996). This name can be traced to an alternative formulation, which tries to find the maximum entropy probability function that satisfies **moment-matching** constraints. These constraints specify that the empirical counts of each label-feature pair should match the expected counts under the induced probability distribution

$$p(y \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}),$$

$$\sum_{i=1}^N f_j(\mathbf{x}^{(i)}, y^{(i)}) = \sum_{i=1}^N \sum_{y \in \mathcal{Y}} p(y \mid \mathbf{x}^{(i)}; \boldsymbol{\theta}) f_j(\mathbf{x}^{(i)}, y), \quad \forall j \quad [1.75]$$

This constraint will be met exactly when the derivative of the conditional log-likelihood function (Equation 1.63) is equal to zero. However, the constraint can be met by many values of $\boldsymbol{\theta}$, so which should we choose?

Let $p_{y|x}$ refer to the conditional distribution *function* of y given \mathbf{x} , and $p(y \mid \mathbf{x})$ to refer to the probability of a specific y given a specific \mathbf{x} . The **entropy** of the conditional probability distribution $p_{y|x}$ is,

$$H(p_{y|x}) = - \sum_{\mathbf{x} \in \mathcal{X}} p_{\mathbf{x}}(\mathbf{x}) \sum_{y \in \mathcal{Y}} p_{y|x}(y \mid \mathbf{x}) \log p_{y|x}(y \mid \mathbf{x}), \quad [1.76]$$

where \mathcal{X} is the set of all possible feature vectors, and $p_{\mathbf{x}}(\mathbf{x})$ is the probability of observing the base features \mathbf{x} . The distribution $p_{\mathbf{x}}$ is unknown, but an empirical estimate can be computed by summing over all the instances in the training set,

$$\tilde{H}(p_{y|x}) = -\frac{1}{N} \sum_{i=1}^N \sum_{y \in \mathcal{Y}} p_{y|x}(y \mid \mathbf{x}^{(i)}) \log p_{y|x}(y \mid \mathbf{x}^{(i)}). \quad [1.77]$$

If the entropy is large, the likelihood function is smooth across possible values of y ; if it is small, the likelihood function is sharply peaked at some preferred value; in the limiting case, the entropy is zero if $p(y \mid \mathbf{x}) = 1$ for some y . The maximum-entropy criterion chooses to make the weakest commitments possible, while satisfying the moment-matching constraints from Equation 1.75. The solution to this constrained optimization problem is identical to the maximum conditional likelihood (logistic-loss) formulation we considered in the previous section.

1.6.3 Further reading

For more on stochastic gradient descent, as applied to a number of different learning algorithms, see (Zhang, 2004) and (Bottou, 1998). Murphy (2012) traces stochastic gradient descent to Nemirovski and Yudin (1978). Kummerfeld et al. (2015) empirically reviews several optimization algorithms for large-margin learning.

1.7 Summary of learning algorithms

Having seen several learning algorithms, it is natural to ask which is best in various situations.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Naïve Bayes *Pros*: easy to implement; estimation is very fast, requiring only a single pass over the data; assigns probabilities to predicted labels; controls overfitting with smoothing parameter. *Cons*: often has poor accuracy, especially with correlated features.

Perceptron *Pros*: easy to implement; online; error-driven learning means that accuracy is typically high, especially after averaging. *Cons*: not probabilistic; non-averaged perceptron performs poorly if data is not separable; hard to know when to stop learning; lack of margin can lead to overfitting.

Support vector machine *Pros*: optimizes an error-based metric, usually resulting in high accuracy; overfitting is controlled by a regularization parameter. *Cons*: not probabilistic.

Logistic regression *Pros*: error-driven and probabilistic; overfitting is controlled by a regularization parameter. *Cons*: batch learning requires black-box optimization; logistic loss sometimes gives lower accuracy than hinge loss, due to overtraining on correctly-labeled examples.

One of the main distinctions is whether the model offers a probability over labels. This is useful in modular architectures, where the output of one classifier is the input for some other system. In cases where probability is not necessary, the support vector machine is usually the right choice, since it is no more difficult to implement than the perceptron, and is often more accurate. When probability is necessary, logistic regression is usually more accurate than Naïve Bayes.

Exercises

1. Derive the maximum-likelihood estimate for the parameter μ in Naïve Bayes.
2. As noted in the discussion of averaged perceptron in § 1.2.2, the computation of the running sum $\mathbf{m} \leftarrow \mathbf{m} + \boldsymbol{\theta}$ is unnecessarily expensive, requiring $K \times V$ operations. Give an alternative way to compute the averaged weights $\bar{\boldsymbol{\theta}}$, with complexity that is independent of V and linear in the sum of feature sizes $\sum_{i=1}^N |\mathbf{f}(\mathbf{x}^{(i)}, y^{(i)})|$.
3. Consider a dataset that is comprised of two identical instances $\mathbf{x}^{(1)} = \mathbf{x}^{(2)}$ with distinct labels $y^{(1)} \neq y^{(2)}$. Assume all features are binary $x_j \in \{0, 1\}$ for all j .
Now suppose that the averaged perceptron always chooses $i = 1$ when t is even, and $i = 2$ when t is odd, and that it will terminate under the following condition:

$$\epsilon \geq \max_j \left| \frac{1}{t} \sum_t \theta_j^{(t)} - \frac{1}{t-1} \sum_t \theta_j^{(t-1)} \right|. \quad [1.78]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

In words, the algorithm stops when the largest change in the averaged weights is less than or equal to ϵ . Compute the number of iterations before the averaged perceptron terminates.

4. Suppose you have two labeled datasets D_1 and D_2 , with the same features and labels.
 - Let $\theta^{(1)}$ be the unregularized logistic regression (LR) coefficients from training on dataset D_1 .
 - Let $\theta^{(2)}$ be the unregularized LR coefficients (same model) from training on dataset D_2 .
 - Let θ^* be the unregularized LR coefficients from training on the combined dataset $D_1 \cup D_2$.

Under these conditions, prove that for any feature j ,

$$\begin{aligned}\theta_j^* &\geq \min(\theta_j^{(1)}, \theta_j^{(2)}) \\ \theta_j^* &\leq \max(\theta_j^{(1)}, \theta_j^{(2)}).\end{aligned}$$

Chapter 2

Nonlinear classification

Linear classification may seem like all we need for natural language processing. The bag-of-words representation is inherently high-dimensional, and the number of features is often larger than the number of training instances. This means that we can usually find a linear classifier that perfectly fits the training data; moving to nonlinear classification may seem to only increase the risk of overfitting. For many tasks, lexical features are meaningful in isolation, and can offer independent evidence about the instance label — unlike computer vision, where individual pixels are rarely informative, and must be evaluated holistically to make sense of an image. For these reasons, natural language processing has historically focused on linear classification to a much greater extent than other machine learning application domains.

Yet in recent years, nonlinear classifiers have swept through natural language processing, and are now default approach for many tasks (Manning, 2016). There are at least three reasons for this change.

- There have been rapid advances in **deep learning**, a family of nonlinear methods that learn complex functions of the input through multiple layers of computation (Goodfellow et al., 2016).
- Deep learning allows facilitates the incorporation of **word embeddings**, which are dense vector representations of words. Word embeddings can be learned from large amounts of unlabeled data, and enable generalization to words that do not appear in the annotated training data (word embeddings are discussed in detail in chapter 13). This makes it possible to train high-coverage NLP systems with relatively little labeled training data.
- A third reason for the rise of deep nonlinear learners is hardware. Standard CPUs are no longer increasing in speed, but many deep learning models can be implemented efficiently on graphics processing units (GPUs), offering substantial perfor-

mance improvements.

This chapter focuses on **neural networks**, which are by far the dominant approach for nonlinear classification in natural language processing today.¹ Historically, a few other nonlinear learning methods have been applied to language data:

- **Kernel methods** are generalizations of the **nearest-neighbor** classification rule, which classifies each instance by the label of the most similar example in the training set (Hastie et al., 2009). The application of the **kernel support vector machine** to information extraction is described in chapter 16.
- **Decision trees** classify instances by checking a set of conditions. Scaling decision trees to bag-of-words inputs is difficult, but they have been effectively applied to problems such as coreference resolution (chapter 14), where more compact feature sets can be constructed (Soon et al., 2001).
- **Boosting** and related **ensemble methods** work by combining the predictions of several “weak” classifiers, each of which may consider only a small subset of features. Boosting has been successfully applied to text classification (Schapire and Singer, 2000) and syntactic analysis (Abney et al., 1999), and remains one of the most competitive methods for machine learning competition sites such as Kaggle (Chen and Guestrin, 2016).

2.1 Feedforward neural networks

Consider the problem of building a classifier for movie reviews. The goal is to predict a label $y \in \{\text{GOOD}, \text{BAD}, \text{OKAY}\}$ from a representation of the text of each document, x . But what makes a good movie? The story, acting, cinematography, soundtrack, and so on. Now suppose the training set contains labels for each of these additional features, $z = [z_1, z_2, \dots, z_{K_z}]^T$. With such information, we could build a two-step classifier:

1. **Use the text x to predict the features z .** Specifically, we can train a logistic regression classifier to compute $p(z_k | x)$, for each $k \in \{1, 2, \dots, K_z\}$.
2. **Use the features z to predict the label y .** Again, we can train a logistic classifier to compute $p(y | z)$. On test data, z is unknown, so we use the probabilities $p(z | x)$ from the first layer as the features.

This setup is shown in Figure 2.1, which describes the proposed classifier in a **computation graph**: the text features x are connected to the middle layer z , which in turn is connected to the label y .

¹I will use “deep learning” and “neural networks” interchangeably.

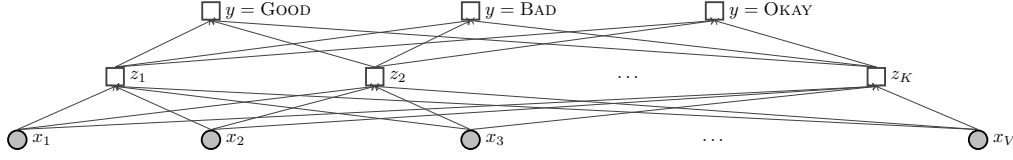


Figure 2.1: A feedforward neural network. Shaded circles indicate observed features, usually words; squares indicate nodes in the computation graph, which are computed from the information carried over the incoming arrows.

Since each $z_k \in \{0, 1\}$, we can treat $p(z_k | \mathbf{x})$ as a binary classification problem, using binary logistic regression:

$$\Pr(z_k = 1 | \mathbf{x}; \Theta^{(x \rightarrow z)}) = \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) = (1 + \exp(-\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}))^{-1}, \quad [2.1]$$

where $\sigma(\cdot)$ is the **sigmoid** function (shown in Figure 2.2), and the matrix $\Theta^{(x \rightarrow z)} \in \mathbb{R}^{K_z \times V}$ is constructed by stacking the weight vectors for each z_k ,

$$\Theta^{(x \rightarrow z)} = [\theta_1^{(x \rightarrow z)}, \theta_2^{(x \rightarrow z)}, \dots, \theta_{K_z}^{(x \rightarrow z)}]^\top. \quad [2.2]$$

We will assume that \mathbf{x} contains a term with a constant value of 1, so that a corresponding offset parameter is included in each $\theta_k^{(x \rightarrow z)}$.

At the output layer, we use the multi-class logistic regression probability,

$$\Pr(y = j | \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \frac{\exp(\theta_j^{(z \rightarrow y)} \cdot \mathbf{z} + b_j)}{\sum_{j' \in \mathcal{Y}} \exp(\theta_{j'}^{(z \rightarrow y)} \cdot \mathbf{z} + b_{j'})}, \quad [2.3]$$

where b_j is an offset for label j , and the output weight matrix $\Theta^{(z \rightarrow y)} \in \mathbb{R}^{K_y \times K_z}$ is again constructed by concatenation,

$$\Theta^{(z \rightarrow y)} = [\theta_1^{(z \rightarrow y)}, \theta_2^{(z \rightarrow y)}, \dots, \theta_{K_y}^{(z \rightarrow y)}]^\top. \quad [2.4]$$

The vector of probabilities over each possible value of y is denoted,

$$\mathbf{p}(y | \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \text{SoftMax}(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b}), \quad [2.5]$$

where element j in the output of the **SoftMax** function is computed as in Equation 2.3.

We have now defined a multilayer classifier, which can be summarized as,

$$\mathbf{p}(z | \mathbf{x}; \Theta^{(x \rightarrow z)}) = \sigma(\Theta^{(x \rightarrow z)} \mathbf{x}) \quad [2.6]$$

$$\mathbf{p}(y | \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \text{SoftMax}(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b}), \quad [2.7]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

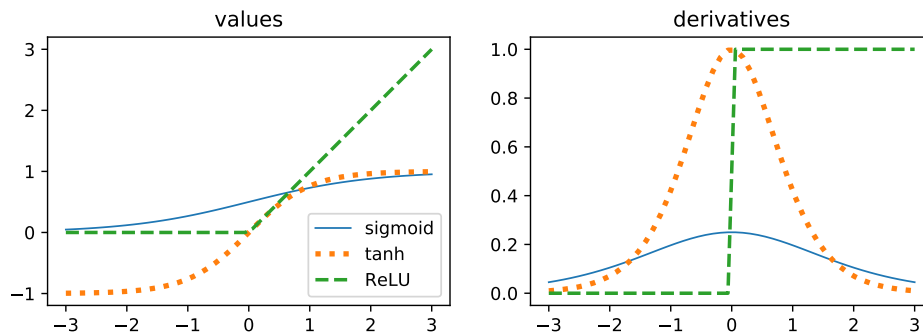


Figure 2.2: The sigmoid, tanh, and ReLU activation functions

where $\sigma(\cdot)$ is now applied **elementwise** to the vector of inner products,

$$\sigma(\Theta^{(x \rightarrow z)} \mathbf{x}) = [\sigma(\theta_1^{(x \rightarrow z)} \cdot \mathbf{x}), \sigma(\theta_2^{(x \rightarrow z)} \cdot \mathbf{x}), \dots, \sigma(\theta_{K_z}^{(x \rightarrow z)} \cdot \mathbf{x})]^\top. \quad [2.8]$$

Now the hidden features \mathbf{z} are never observed, even in the training data. We can still construct the architecture in Figure 2.1. Instead of predicting y from a discrete vector of predicted values \mathbf{z} , we use the probabilities $\sigma(\theta_k \cdot \mathbf{x})$. The resulting classifier is barely changed:

$$\mathbf{z} = \sigma(\Theta^{(x \rightarrow z)} \mathbf{x}) \quad [2.9]$$

$$p(y \mid \mathbf{x}; \Theta^{(z \rightarrow y)}, \mathbf{b}) = \text{SoftMax}(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b}). \quad [2.10]$$

We have now defined a classification model that predicts the label $y \in \mathcal{Y}$ from the base features \mathbf{x} , through a “hidden layer” \mathbf{z} . This is a **feedforward neural network**.²

2.2 Neural networks

This feedforward neural network can be generalized in a number of ways.

2.2.1 Activation functions

If the hidden layer is viewed as a set of latent features, then the sigmoid function represents the extent to which each of these features is “activated” by a given input. However, we can relax this view of the hidden layer, and regard it more generally as a nonlinear transformation of the input. This opens the door to many other activation functions, some of which are shown in Figure 2.2. At the moment, the choice of activation functions is more art than science, but a few points can be made about the most popular varieties:

²The architecture is sometimes called a **multilayer perceptron**, but this is misleading, because each layer is not a perceptron as defined in Algorithm 3.

- The range of the sigmoid function is $(0, 1)$. The bounded range ensures that a cascade of sigmoid functions will not “blow up” to a huge output, and this is important for deep networks with several hidden layers. The derivative of the sigmoid is $\frac{\partial}{\partial a} \sigma(a) = \sigma(a)(1 - \sigma(a))$. This derivative becomes small at the extremes, which can make learning slow; this is called the **vanishing gradient** problem.
- The range of the **tanh activation function** is $(-1, 1)$: like the sigmoid, the range is bounded, but unlike the sigmoid, it includes negative values. The derivative is $\frac{\partial}{\partial a} \tanh(a) = 1 - \tanh(a)^2$, which is steeper than the logistic function near the origin (LeCun et al., 1998). The tanh function can also suffer from vanishing gradients at extreme values.
- The **rectified linear unit (ReLU)** is zero for negative inputs, and linear for positive inputs (Glorot et al., 2011),

$$\text{ReLU}(a) = \begin{cases} a, & a \geq 0 \\ 0, & \text{otherwise.} \end{cases} \quad [2.11]$$

The derivative is a step function, which is 1 if the input is positive, and zero otherwise. Once the activation is zero, the gradient is also zero. This can lead to the problem of **dead neurons**, where some ReLU nodes are zero for all inputs, throughout learning. A solution is the **leaky ReLU**, which has a small positive slope for negative inputs (Maas et al., 2013),

$$\text{Leaky-ReLU}(a) = \begin{cases} a, & a \geq 0 \\ .0001a, & \text{otherwise.} \end{cases} \quad [2.12]$$

Sigmoid and tanh are sometimes described as **squashing functions**, because they squash an unbounded input into a bounded range. Glorot and Bengio (2010) recommend against the use of the sigmoid activation in deep networks, because its mean value of $\frac{1}{2}$ can cause the next layer of the network to be saturated, with very small gradients on their own parameters. Several other activation functions are reviewed by Goodfellow et al. (2016), who recommend ReLU as the “default option.”

2.2.2 Network structure

Deep networks stack up several hidden layers, $z^{(1)}, z^{(2)}, \dots$, where each $z^{(d)}$ is the input to the next layer, $z^{(d+1)}$. As the total number of nodes in the network increases, so does its capacity to learn complex functions of the input. For a fixed number of nodes, an architectural decision is whether to emphasize width (large K_z at each layer) or depth (many layers). At present, this tradeoff is not well understood.³

³There is a universal approximation theorem that states that with even a single hidden layer, a neural network can approximate any continuous function on a closed and bounded subset of \mathbb{R}^N to an arbitrarily

It is also possible to “short circuit” a hidden layer, by propagating information directly from the input to the next higher level of the network. This is the idea behind **residual networks**, which propagate information directly from the input to the subsequent layer (He et al., 2016),

$$\mathbf{z} = f(\Theta^{(x \rightarrow z)} \mathbf{x}) + \mathbf{x}, \quad [2.13]$$

where f is any nonlinearity, such as sigmoid or ReLU. A more complex architecture is the **highway network** (Srivastava et al., 2015; Kim et al., 2016), in which an addition **gate** controls an interpolation between $f(\Theta^{(x \rightarrow z)} \mathbf{x})$ and \mathbf{x} :

$$\mathbf{t} = \sigma(\Theta^{(t)} \mathbf{x} + \mathbf{b}^{(t)}) \quad [2.14]$$

$$\mathbf{z} = \mathbf{t} \odot f(\Theta^{(x \rightarrow z)} \mathbf{x}) + (\mathbf{1} - \mathbf{t}) \odot \mathbf{x}, \quad [2.15]$$

where \odot refers to an elementwise vector product, and $\mathbf{1}$ is a column vector of ones. Gating is also used in the **long short-term memory (LSTM)**, which is discussed in chapter 5. Residual and highway connections address a problem with a deep architectures: repeated application of a nonlinear activation function can make it difficult to learn the parameters of the lower levels of the network, which are too distant from the supervision signal.

2.2.3 Outputs and loss functions

In the multi-class classification example, a softmax output produces probabilities over each possible label. This aligns with a negative **conditional log-likelihood** loss function,

$$-\mathcal{L} = - \sum_{i=1}^N \log p(y^{(i)} | \mathbf{x}^{(i)}; \Theta). \quad [2.16]$$

where $\Theta = \{\Theta^{(x \rightarrow z)}, \Theta^{(z \rightarrow y)}, \mathbf{b}\}$ is the entire set of parameters.

This loss can be written alternatively as follows:

$$\tilde{y}_j \triangleq \Pr(y = j | \mathbf{x}^{(i)}; \Theta) \quad [2.17]$$

$$-\mathcal{L} = - \sum_{i=1}^N \mathbf{e}_{y^{(i)}} \cdot \log \tilde{\mathbf{y}} \quad [2.18]$$

where $\mathbf{e}_{y^{(i)}}$ is a “one-hot” vector of zeros with a value of 1 at position $y^{(i)}$. The inner product between $\mathbf{e}_{y^{(i)}}$ and $\log \tilde{\mathbf{y}}$ is also called the multinomial **cross-entropy**, and this terminology is preferred in many neural networks papers and software packages.

small non-zero error; see section 6.4.1 of Goodfellow et al. (2016) for a survey of these theoretical results. However, the width of the hidden layer may need to be arbitrarily large. Furthermore, the fact that a network has the capacity to approximate any given function does not say anything about whether it is possible to *learn* the function using gradient-based optimization.

(c) Jacob Eisenstein 2014-2017. Work in progress.

It is also possible to train neural networks from other objectives, such as a margin loss. In this case, it is not necessary to use softmax at the output layer: an affine transformation of the hidden layer is enough:

$$\Psi(y; \mathbf{x}^{(i)}, \Theta) = \boldsymbol{\theta}_y^{(z \rightarrow y)} \cdot \mathbf{z} + b_y \quad [2.19]$$

$$\ell_{\text{MARGIN}}(\Theta; \mathbf{x}^{(i)}, y^{(i)}) = \max_{y \neq y^{(i)}} \left(1 + \Psi(y; \mathbf{x}^{(i)}, \Theta) - \Psi(y^{(i)}; \mathbf{x}^{(i)}, \Theta) \right). \quad [2.20]$$

In regression problems, the output is a scalar or vector (see § 3.1.2). For these problems, a typical loss function is the squared error $(y - \hat{y})^2$ or squared norm $\|\mathbf{y} - \hat{\mathbf{y}}\|_2^2$.

2.2.4 Inputs and lookup layers

In text classification, the input layer \mathbf{x} can refer to a bag-of-words vector, where x_j is the count of word j . The input to the hidden unit z_k is then $\sum_{j=1}^V \theta_{j,k}^{(x \rightarrow z)} x_j$, and word j is represented by the vector $\boldsymbol{\theta}_j^{(x \rightarrow z)}$. This vector is sometimes described as the **embedding** of word j , and can be learned from unlabeled data, using techniques discussed in chapter 13. The columns of $\Theta^{(x \rightarrow z)}$ are each K_z -dimensional word embeddings.

In chapter 1, we encountered another view of text documents, as a sequence of word tokens, w_1, w_2, \dots, w_M . To implement this view in a neural network, represent each word token w_m with a one-hot vector, $\mathbf{e}_{w_m} \in \mathbb{R}^V$. The matrix-vector product $\Theta^{(x \rightarrow z)} \mathbf{e}_{w_m}$ returns the embedding of word w_m . The complete document can be represented by horizontally concatenating these one-hot vectors, $\mathbf{W} = [\mathbf{e}_{w_1}, \mathbf{e}_{w_2}, \dots, \mathbf{e}_{w_M}]$, and the bag-of-words representation can be recovered from the matrix-vector product $\mathbf{W}\mathbf{1}$, which simply sums over the tokens $m = \{1, 2, \dots, M\}$. The matrix product $\Theta^{(x \rightarrow z)} \mathbf{W}$ contains the horizontally concatenated embeddings of each word in the document, which will be useful the starting point for **convolutional neural networks** (§ 2.4). This is sometimes called a **lookup layer**, because the first step is to lookup the embeddings for each word in the input text.

2.3 Learning neural networks

The feedforward network in Figure 2.1 can now be written in a more general form,

$$\mathbf{z} \leftarrow f(\Theta^{(x \rightarrow z)} \mathbf{x}^{(i)}) \quad [2.21]$$

$$\tilde{\mathbf{y}} \leftarrow \text{SoftMax} \left(\Theta^{(z \rightarrow y)} \mathbf{z} + \mathbf{b} \right) \quad [2.22]$$

$$\ell^{(i)} \leftarrow -\mathbf{e}_{y^{(i)}} \cdot \log \tilde{\mathbf{y}}, \quad [2.23]$$

where f is an elementwise activation function, such as σ or ReLU.

(c) Jacob Eisenstein 2014-2017. Work in progress.

We now consider how to estimate the parameters: $\Theta^{(x \rightarrow z)}$, $\Theta^{(z \rightarrow y)}$ and \mathbf{b} . As in chapter 1, we will focus on online gradient-based optimization to incrementally minimize a local loss function with respect to each parameter. The simplest such algorithm is stochastic gradient descent,

$$\mathbf{b} \leftarrow \mathbf{b} - \eta^{(t)} \nabla_{\mathbf{b}} \ell^{(i)} \quad [2.24]$$

$$\theta_k^{(z \rightarrow y)} \leftarrow \theta_k^{(z \rightarrow y)} - \eta^{(t)} \nabla_{\theta_k^{(z \rightarrow y)}} \ell^{(i)} \quad [2.25]$$

$$\theta_k^{(x \rightarrow z)} \leftarrow \theta_k^{(x \rightarrow z)} - \eta^{(t)} \nabla_{\theta_k^{(x \rightarrow z)}} \ell^{(i)}, \quad [2.26]$$

where $\eta^{(t)}$ is the learning rate on iteration t , and $\ell^{(i)}$ is the loss at instance (or minibatch) i .

The gradients of the negative log-likelihood on \mathbf{b} and $\theta_k^{(z \rightarrow y)}$ are very similar to the gradients in logistic regression,

$$\nabla_{\theta_k^{(z \rightarrow y)}} \ell^{(i)} = \left[\frac{\partial \ell^{(i)}}{\partial \theta_{k,1}^{(z \rightarrow y)}}, \frac{\partial \ell^{(i)}}{\partial \theta_{k,2}^{(z \rightarrow y)}}, \dots, \frac{\partial \ell^{(i)}}{\partial \theta_{k,K_y}^{(z \rightarrow y)}} \right]^\top \quad [2.27]$$

$$\frac{\partial \ell^{(i)}}{\partial \theta_{k,j}^{(z \rightarrow y)}} = - \frac{\partial}{\partial \theta_{k,j}^{(z \rightarrow y)}} \left(\theta_{y^{(i)}}^{(z \rightarrow y)} \cdot \mathbf{z} - \log \sum_{y \in \mathcal{Y}} \exp \theta_y^{(z \rightarrow y)} \cdot \mathbf{z} \right) \quad [2.28]$$

$$= \left(\Pr(y = j \mid \mathbf{z}; \Theta^{(z \rightarrow y)}, \mathbf{b}) - \delta(j = y^{(i)}) \right) z_k, \quad [2.29]$$

where $\delta(j = y^{(i)})$ is a function that returns one when $j = y^{(i)}$, and zero otherwise. The gradient $\nabla_{\mathbf{b}} \ell^{(i)}$ is similar to Equation 2.29.

The gradients on the input layer weights $\Theta^{(x \rightarrow z)}$ can be obtained by applying the chain rule of differentiation:

$$\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} = \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial z_k}{\partial \theta_{n,k}^{(x \rightarrow z)}} \quad [2.30]$$

$$= \frac{\partial \ell^{(i)}}{\partial z_k} \frac{\partial f(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})}{\partial \theta_{n,k}^{(x \rightarrow z)}} \quad [2.31]$$

$$= \left(\frac{\partial \ell^{(i)}}{\partial z_k} \right) \times f'(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times x_n, \quad [2.32]$$

where $f'(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})$ is the derivative of the activation function f , applied at the input

$\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}$. For example, if f is the sigmoid function, then the derivative is,

$$\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} = \left(\frac{\partial \ell^{(i)}}{\partial z_k} \right) \times \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x}) \times (1 - \sigma(\theta_k^{(x \rightarrow z)} \cdot \mathbf{x})) \times x_n \quad [2.33]$$

$$= \left(\frac{\partial \ell^{(i)}}{\partial z_k} \right) \times z_k \times (1 - z_k) \times x_n. \quad [2.34]$$

For intuition, consider each of the terms in the product.

- If the negative log-likelihood $\ell^{(i)}$ does not depend much on z_k , $\frac{\partial \ell^{(i)}}{\partial z_k} \rightarrow 0$, then it doesn't matter how we compute z_k , and so $\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} \rightarrow 0$.
- If z_k is near 1 or 0, then the curve of the sigmoid function (Figure 2.2) is nearly flat, and changing the inputs will make little local difference. The term $z_k \times (1 - z_k)$ is maximized at $z_k = \frac{1}{2}$, where the slope of the sigmoid function is steepest.
- If $x_n = 0$, then it does not matter how we set the weights $\theta_{n,k}^{(x \rightarrow z)}$, so $\frac{\partial \ell^{(i)}}{\partial \theta_{n,k}^{(x \rightarrow z)}} = 0$.

2.3.1 Backpropagation

In the equations above, the value $\frac{\partial \ell^{(i)}}{\partial z_k}$ is reused in the derivatives with respect to each $\theta_{n,k}^{(x \rightarrow z)}$. It should therefore be computed once, and then cached. Furthermore, we should only compute any derivative once we have already computed all of the necessary “inputs” demanded by the chain rule of differentiation. This combination of sequencing, caching, and differentiation is known as **backpropagation**. It can be generalized to any directed acyclic **computation graph**.

A computation graph is a declarative representation of a computational process. At each node t , compute a value v_t by applying a function f_t to a (possibly empty) list of parent nodes, π_t . For example, in a feedforward network with one hidden layer, there are nodes for the inputs $\mathbf{x}^{(i)}$, the hidden layer \mathbf{z} , the predicted output $\tilde{\mathbf{y}}$, and the parameters $\{\Theta^{(x \rightarrow z)}, \Theta^{(z \rightarrow y)}, \mathbf{b}\}$. During training, there is also a node for the observed label $\mathbf{y}^{(i)}$ and the loss $\ell^{(i)}$. Various types of nodes play different roles during learning:

- *Variables* include the *inputs* \mathbf{x} , the *hidden nodes* \mathbf{z} , the *outputs* \mathbf{y} , and the loss function. Inputs are variables that do not have parents. Backpropagation computes the gradients with respect to all variables except the inputs, but does not update the variables during learning.
- *Parameters* include the weights and offsets. They do not have parents, and they are updated during learning.

Algorithm 6 General backpropagation algorithm. In the computation graph G , every node contains a function f_t and a set of inputs π_t ; the inputs to the graph are $\mathbf{x}^{(i)}$. [todo: needs work]

```

1: procedure BACKPROP( $G = \{f_t, \pi_t\}_{t=1}^T, \mathbf{x}^{(i)}$ )
2:    $v_{t(n)} \leftarrow x_n^{(i)}$  for all  $n$  and associated computation nodes  $t(n)$ .
3:   for  $t \in \text{TOPOLOGICALSORT}(G)$  do    ▷ Forward pass: compute value at each node
4:     if  $|\pi_t| > 0$  then
5:        $v_t \leftarrow f_t(v_{\pi_{t,1}}, v_{\pi_{t,2}}, \dots, v_{\pi_{t,N_t}})$ 
6:    $g_{\text{objective}} = 1$     ▷ Backward pass: compute gradients at each node
7:   for  $t \in \text{REVERSE}(\text{TOPOLOGICALSORT}(G))$  do
8:      $g_t \leftarrow \sum_{t': t \in \pi_{t'}} g_{t'} \frac{\partial f_{t'}(v_{\pi_{t',1}}, v_{\pi_{t',2}}, \dots, v_{\pi_{t',N_{t'}}})}{\partial v_t}$ 
9:   return  $\{g_1, g_2, \dots, g_T\}$ 

```

- The *objective* is not the parent of any other node. Backpropagation begins by computing the gradient with respect to this node.

If the computation graph is a directed acyclic graph, then it is possible to order the nodes with a topological sort, so that if node t is an input to node t' , then $t < t'$. This means that the values $\{v_t\}_{t=1}^T$ can be computed in a single forward pass. The reverse of the topological sort is used to compute gradients: each gradient g_t is computed from the gradients of nodes to which t is an input, implementing the chain rule of differentiation. The general backpropagation algorithm for computation graphs is shown in Algorithm 6, and illustrated in Figure 2.3.

While the gradients with respect to each parameter may be complex, they are composed of products of simple parts. For many networks, all gradients can be computed through **automatic differentiation**. This means that end users need only specify the feed-forward computation, and the gradients necessary for learning can be obtained automatically. There are many software libraries that perform automatic differentiation on computation graphs, such as `Torch` (Collobert et al., 2011a), `TensorFlow` (Abadi et al., 2016), and `DyNet` (Neubig et al., 2017a). One important distinction between these libraries is whether they support **dynamic computation graphs**, in which the structure of the computation graph varies across instances. Static computation graphs are compiled in advance, and can be applied to fixed-dimensional data, such as bag-of-words vectors. In many natural language processing problems, it is desirable to represent the input as a sequence or even a tree structure, resulting in a variable computation graph at each instance.

(c) Jacob Eisenstein 2014-2017. Work in progress.

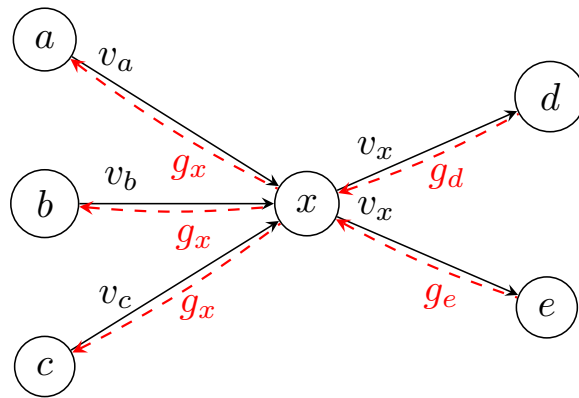


Figure 2.3: Backpropagation at a single node x in the computation graph. The values of the predecessors v_a, v_b, v_c are the inputs to x , which computes v_x , and passes it on to the successors d and e . The gradients at the successors g_d and g_e are passed back to x , where they are incorporated into the gradient g_x , which is then passed back to the predecessors a, b , and c .

2.3.2 Regularization and dropout

In linear classification, overfitting was addressed by augmenting the objective with a regularization term, $\lambda \|\theta\|_2^2$. This same approach can be applied to feedforward neural networks, penalizing each matrix of weights:

$$L = \sum_{i=1}^N \ell^{(i)} + \lambda_{z \rightarrow y} \|\Theta^{(z \rightarrow y)}\|_F^2 + \lambda_{x \rightarrow z} \|\Theta^{(x \rightarrow z)}\|_F^2, \quad [2.35]$$

where $\|\Theta\|_F^2 = \sum_{i,j} \theta_{i,j}^2$ is the squared **Frobenius norm**, which generalizes the L_2 norm to matrices. The bias parameters \mathbf{b} are not regularized, as they do not contribute to the sensitivity of the classifier to the inputs. In gradient-based optimization, the practical effect of Frobenius norm regularization is that the weights “decay” towards zero at each update, motivating the alternative name **weight decay**.

Another approach to controlling model complexity is **dropout**, which involves randomly setting some computation nodes to zero during training (Srivastava et al., 2014). For example, in the feedforward network, on each training instance, with probability ρ we set each input x_n and each hidden layer node z_k to zero. Srivastava et al. (2014) recommend $\rho = 0.5$ for hidden units, and $\rho = 0.2$ for input units. Dropout is also incorporated in the gradient computation, so if node z_k is dropped, then none of the weights $\theta_k^{(x \rightarrow z)}$ will be updated for this instance. Dropout prevents the network from learning to depend too much on any one feature or hidden node, and prevents **feature co-adaptation**, in which a hidden unit is only useful in combination with one or more other hidden units. Dropout is

a special case of **feature noising**, which can also involve adding Gaussian noise to inputs or hidden units (Holmstrom and Koistinen, 1992). Wager et al. (2013) show that dropout is approximately equivalent to “adaptive” L_2 regularization, with a separate regularization penalty for each feature.

2.3.3 *Learning theory

Chapter 1 emphasized the importance of **convexity** for learning: for convex objectives, the global optimum can be found efficiently. The negative log-likelihood and hinge loss are convex functions of the parameters of the output layer. However, the output of a feed-forward network is generally not a convex function of the parameters of the input layer, $\Theta^{(x \rightarrow z)}$. Feedforward networks can be viewed as function composition, where each layer is a function that is applied to the output of the previous layer. Convexity is generally not preserved in the composition of two convex functions — and furthermore, “squashing” activation functions like tanh and sigmoid are not convex.

Here is a more intuitive explanation for why the objective of a neural network with one or more hidden layers is not a convex function of its inputs. Apply a permutation π to the elements in the hidden layer, from $\mathbf{z} = [z_1, z_2, \dots, z_{K_z}]$ to $\tilde{\mathbf{z}} = [z_{\pi(1)}, z_{\pi(2)}, \dots, z_{\pi(K_z)}]$. This corresponds to applying π to the rows of $\Theta^{(x \rightarrow z)}$ and the columns of $\Theta^{(z \rightarrow y)}$, resulting in permuted parameter matrices $\Theta_\pi^{(x \rightarrow z)}$ and $\Theta_\pi^{(z \rightarrow y)}$. As long as this permutation is applied consistently, the loss will be identical, $L(\Theta) = L(\Theta_\pi)$: it is *invariant* to this permutation. However, the loss of the linear combination $L(\alpha\Theta + (1 - \alpha)\Theta_\pi)$ will generally not be identical to the loss under Θ or its permutations. If $L(\Theta)$ is better than the loss at any points in the immediate vicinity, and if $L(\Theta) = L(\Theta_\pi)$, then the loss function does not satisfy the definition of convexity (see § 1.3). One of the exercises asks you to prove this more rigorously.

In practice, the existence of multiple optima is not necessary problematic, if all such optima are permutations of the sort described in the previous paragraph. In contrast, “bad” local optima are better than their neighbors, but much worse than the global optimum. Fortunately, in large feedforward neural networks, most local optima are nearly as good as the global optimum (Choromanska et al., 2015), which helps to explain why backpropagation works. More generally, a **critical point** is one at which the gradient is zero (this is also called a **stationary point**). Critical points may be local optima, but they may also be **saddle points**, which are local minima in some directions, but local *maxima* in other directions. For example, the equation $x_1^2 - x_2^2$ has a saddle point at $\mathbf{x} = (0, 0)$.⁴ In large networks, the overwhelming majority of critical points are saddle points, rather than local minima or maxima (Dauphin et al., 2014). Saddle points can pose problems for gradient-based optimization, since learning will slow to a crawl as the gradient goes

⁴Thanks to Rong Ge’s blogpost for this example, <http://www.offconvex.org/2016/03/22/saddlepoints/>

to zero. However, the noise introduced by stochastic gradient descent (and also by feature noising techniques such as dropout) can help online optimization to escape saddle points and find high-quality optima (Ge et al., 2015). Other techniques address saddle points directly, using local reconstructions of the Hessian matrix (Dauphin et al., 2014) or higher-order derivatives (Anandkumar and Ge, 2016).

2.3.4 Tricks

Getting neural networks to work effectively sometimes requires heuristic “tricks” (Bottou, 2012; Goodfellow et al., 2016; Goldberg, 2017b). This section presents some tricks that are especially important.

Initialization Initialization is not especially important for linear classifiers, since convexity ensures that the global optimum can usually be found quickly. But for multilayer neural networks, it is helpful to have a good starting point. One reason is that if the magnitude of the initial weights is too large, a sigmoid or tanh nonlinearity will be saturated, leading to a small gradient, and slow learning. Large gradients are also problematic. Initialization can help avoid these problems, by ensuring that the variance over the initial gradients is constant and bounded throughout the network. For networks with tanh activation functions, this can be achieved by sampling the initial weights from the following uniform distribution (Glorot and Bengio, 2010),

$$\theta_{i,j} \sim U \left[-\frac{\sqrt{6}}{\sqrt{d_{\text{in}}(n) + d_{\text{out}}(n)}}, \frac{\sqrt{6}}{\sqrt{d_{\text{in}}(n) + d_{\text{out}}(n)}} \right], \quad [2.36]$$

[2.37]

For the weights leading to a ReLU activation function, He et al. (2015) use similar argumentation to justify sampling from a zero-mean Gaussian distribution,

$$\theta_{i,j} \sim N(0, \sqrt{2/d_{\text{in}}(n)}) \quad [2.38]$$

Rather than initializing the weights independently, it can be beneficial to initialize each layer jointly as an **orthonormal matrix**, ensuring that $\Theta^\top \Theta = \mathbf{I}$ (Saxe et al., 2014). Orthonormal matrices preserve the norm of the input, so that $\|\Theta \mathbf{x}\| = \|\mathbf{x}\|$, which prevents the gradients from exploding or vanishing. Orthogonality ensures that the hidden units are uncorrelated, so that they correspond to different features of the input. Orthonormal initialization can be performed by applying **singular value decomposition** to a matrix of

(c) Jacob Eisenstein 2014-2017. Work in progress.

values sampled from a standard normal distribution:

$$a_{i,j} \sim N(0, 1) \quad [2.39]$$

$$\mathbf{A} = \{a_{i,j}\}_{i=1,j=1}^{d_{\text{in}}(j), d_{\text{out}}(j)} \quad [2.40]$$

$$\mathbf{U}, \mathbf{S}, \mathbf{V}^\top = \text{SVD}(\mathbf{A}) \quad [2.41]$$

$$\Theta^{(j)} \leftarrow \mathbf{U}. \quad [2.42]$$

The matrix \mathbf{U} contains the **singular vectors** of \mathbf{A} , and is guaranteed to be orthonormal. For more on singular value decomposition, see chapter 13.

Even with careful initialization, there can still be significant variance in the final results. It can be useful to make multiple training runs, and select the one with the best performance on a heldout development set.

Clipping and normalizing the gradients As already discussed, the magnitude of the gradient can pose problems for learning: too large, and learning can diverge, with successive updates thrashing between increasingly extreme values; too small, and learning can grind to a halt. Several heuristics have been proposed to address this issue.

- In **gradient clipping** (Pascanu et al., 2013), an upper limit is placed on the norm of the gradient, and the gradient is rescaled when this limit is exceeded,

$$\text{CLIP}(\tilde{\mathbf{g}}) = \begin{cases} \mathbf{g} & \|\hat{\mathbf{g}}\| < \tau \\ \frac{\tau}{\|\mathbf{g}\|} \mathbf{g} & \text{otherwise.} \end{cases} \quad [2.43]$$

- In **batch normalization** (Ioffe and Szegedy, 2015), the inputs to each computation node are recentered by their mean and variance across all of the instances in the minibatch \mathcal{B} (see § 1.5.2). For example, in a feedforward network with one hidden layer, batch normalization would transform the inputs to the hidden layer as follows:

$$\boldsymbol{\mu}^{(\mathcal{B})} = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} \mathbf{x}^{(i)} \quad [2.44]$$

$$\mathbf{s}^{(\mathcal{B})} = \frac{1}{|\mathcal{B}|} \sum_{i \in \mathcal{B}} (\mathbf{x}^{(i)} - \boldsymbol{\mu}^{(\mathcal{B})})^2 \quad [2.45]$$

$$\bar{\mathbf{x}}^{(i)} = (\mathbf{x}^{(i)} - \boldsymbol{\mu}^{(\mathcal{B})}) / \sqrt{\mathbf{s}^{(\mathcal{B})}}. \quad [2.46]$$

Empirically, this speeds convergence of deep architectures. One explanation is that it helps to correct for changes in the distribution of activations during training.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- In **layer normalization** (Ba et al., 2016), the inputs to each nonlinear activation function are recentered across the layer:

$$\mathbf{a} = \Theta^{(x \rightarrow z)} \mathbf{x} \quad [2.47]$$

$$\mu = \frac{1}{K_z} \sum_{k=1}^{K_z} a_k \quad [2.48]$$

$$s = \frac{1}{K_z} \sum_{k=1}^{K_z} (a_k - \mu)^2 \quad [2.49]$$

$$\mathbf{z} = (\mathbf{a} - \mu) / \sqrt{s}. \quad [2.50]$$

Layer normalization has similar motivations to batch normalization, but it can be applied across a wider range of architectures and training conditions.

Online optimization The trend towards deep learning has spawned a cottage industry of online optimization algorithms, which attempt to improve on stochastic gradient descent. **AdaGrad** was reviewed in § 1.5.2; its main innovation is to set adaptive learning rates for each parameter by storing the sum of squared gradients. Rather than using the sum over the entire training history, we can keep a running estimate,

$$v_j^{(t)} = \beta v_j^{(t-1)} + (1 - \beta) g_{t,j}^2, \quad [2.51]$$

where $g_{t,j}$ is the gradient with respect to parameter j at time t , and $\beta \in [0, 1]$. This term places more emphasis on recent gradients, and is employed in the **AdaDelta** (Zeiler, 2012) and **Adam** (Kingma and Ba, 2014) optimizers. Online optimization and its theoretical background are reviewed by Bottou et al. (2016). **Early stopping**, mentioned in § 1.2.2, can help to avoid overfitting, by terminating training after reaching a plateau in the performance on a heldout validation set.

2.4 Convolutional neural networks

A basic weakness of the bag-of-words model is its inability to account for the ways in which words combine to create meaning, including even simple reversals such as *not pleasant*, *hardly a generous offer*, and *I wouldn't mind missing the flight*. Similarly, computer vision faces the challenge of identifying the semantics of images from pixel features that are profoundly uninformative in isolation. An earlier generation of computer vision research focused on designing *filters* to aggregate local pixel-level features into more meaningful representations, such as edges and corners (e.g., Canny, 1987). Similarly, earlier NLP research attempted to capture multiword linguistic phenomena by hand-designed lexical patterns (Hobbs et al., 1997). In both cases, the output of the filters and patterns

(c) Jacob Eisenstein 2014-2017. Work in progress.

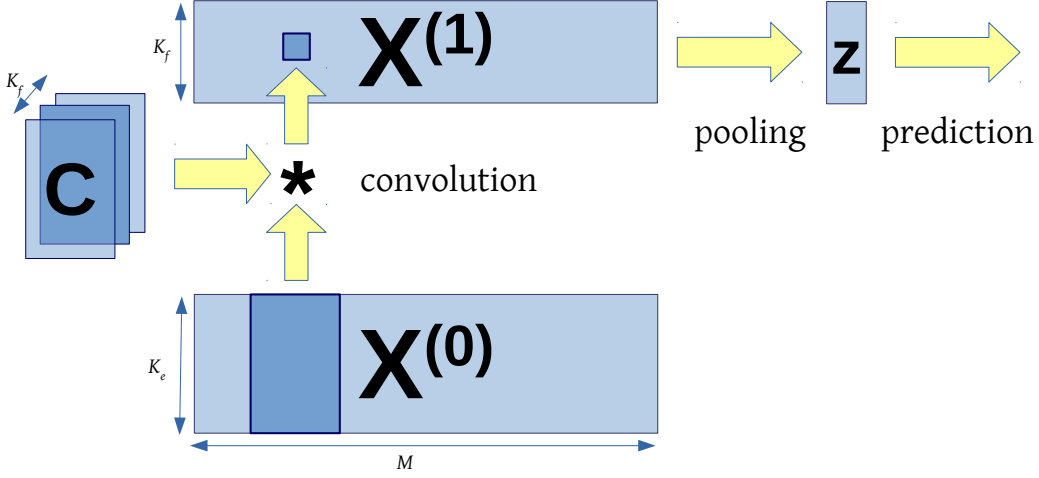


Figure 2.4: A convolutional neural network for text classification

could then act as base features in a linear classifier. But rather than designing these feature extractors by hand, a better approach is to learn them, using the magic of backpropagation. This is the idea behind **convolutional neural networks**.

Following § 2.2.4, define the base layer of a neural network as $\mathbf{X}^{(0)} = \Theta^{(x \rightarrow z)}[e_{w_1}, e_{w_2}, \dots, e_{w_M}]$, where e_{w_m} is a column vector of zeros, with a 1 at position w_m . The base layer has dimension $\mathbf{X}^{(0)} \in \mathbb{R}^{K_e \times M}$, where K_e is the size of the word embeddings. To merge information across adjacent words, we *convolve* $\mathbf{X}^{(0)}$ with a set of filter matrices $\mathbf{C}^{(k)} \in \mathbb{R}^{K_f \times h}$. Convolution is indicated by the symbol $*$, and is defined,

$$\mathbf{X}^{(1)} = f(\mathbf{b} + \mathbf{C} * \mathbf{X}^{(0)}) \implies x_{k,m}^{(1)} = f\left(b_k + \sum_{k'=1}^{K_e} \sum_{n=1}^h c_{k',n}^{(k)} \times x_{k',m+n-1}^{(0)}\right), \quad [2.52]$$

where f is an activation function such as tanh or ReLU, and \mathbf{b} is a vector of offsets. The convolution operation slides the matrix $\mathbf{C}^{(k)}$ across the columns of $\mathbf{X}^{(0)}$; at each position m , compute the elementwise product $\mathbf{C}^{(k)} \odot \mathbf{X}_{m:m+h-1}^{(0)}$, and take the sum.

A simple filter might compute a weighted average over nearby words,

$$\mathbf{C}^{(k)} = \begin{bmatrix} 0.5 & 1 & 0.5 \\ 0.5 & 1 & 0.5 \\ \dots & \dots & \dots \\ 0.5 & 1 & 0.5 \end{bmatrix}, \quad [2.53]$$

thereby representing trigram units like *not so unpleasant*. In **one-dimensional convolution**, each filter matrix $\mathbf{C}^{(k)}$ is constrained to have non-zero values only at row k (Kalchbrenner et al., 2014).

(c) Jacob Eisenstein 2014-2017. Work in progress.

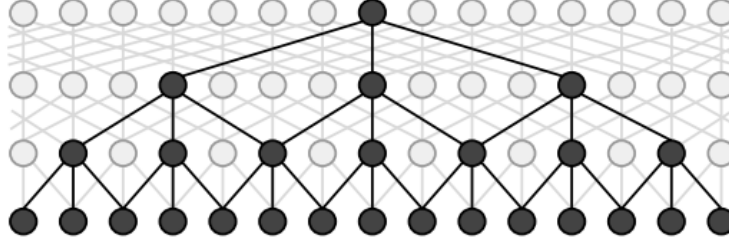


Figure 2.5: A dilated convolutional neural network captures progressively larger context through recursive application of the convolutional operator (Strubell et al., 2017) [todo: permission]

The base matrix $\mathbf{X}^{(0)}$ may be padded with h column vectors of zeros at the beginning and end; this is known as **wide convolution**. If padding is not applied, then the output from each layer will be $h - 1$ units smaller than the input; this is known as **narrow convolution**. The filter matrices need not have identical filter widths, so more generally we could write h_k to indicate the width of filter $\mathbf{C}^{(k)}$. As suggested by the notation $\mathbf{X}^{(0)}$, multiple layers of convolution may be applied, so that $\mathbf{X}^{(d)}$ is the input to $\mathbf{X}^{(d+1)}$.

After D convolutional layers, we obtain a matrix representation of the document $\mathbf{X}^{(D)} \in \mathbb{R}^{K_z \times M}$. If the instances have variable lengths, it is necessary to aggregate over all M word positions to obtain a fixed-length representation. This can be done by a **pooling** operation, such as max-pooling (Collobert et al., 2011b) or average-pooling,

$$\mathbf{z} = \text{MaxPool}(\mathbf{X}^{(D)}) \implies z_k = \max(x_{k,1}^{(D)}, x_{k,2}^{(D)}, \dots, x_{k,M}^{(D)}) \quad [2.54]$$

$$\mathbf{z} = \text{AvgPool}(\mathbf{X}^{(D)}) \implies z_k = \frac{1}{M} \sum_{m=1}^M x_{k,m}^{(D)}. \quad [2.55]$$

The vector \mathbf{z} can now act as a layer in a feedforward network, culminating in a prediction \hat{y} and a loss $\ell^{(i)}$. The setup is shown in Figure 2.4.

Just as in feedforward networks, the parameters $(\mathbf{C}^{(k)}, \mathbf{b}, \Theta)$ can be learned by backpropagating from the classification loss. This requires backpropagating through the max-pooling operation, which is a discontinuous function of the input. But because we need only a local gradient, backpropagation flows only through the argmax m :

$$\frac{\partial z_k}{\partial x_{k,m}^{(D)}} = \begin{cases} 1, & x_{k,m}^{(D)} = \max(x_{k,1}^{(D)}, x_{k,2}^{(D)}, \dots, x_{k,M}^{(D)}) \\ 0, & \text{otherwise.} \end{cases} \quad [2.56]$$

The computer vision literature has produced a huge variety of convolutional architectures, and many of these bells and whistles can be applied to text data. One avenue for

improvement is more complex pooling operations, such as k -max pooling (Kalchbrenner et al., 2014), which returns a matrix of the k largest values for each filter. Another innovation is the use of **dilated convolution** to build multiscale representations (Yu and Koltun, 2016). At each layer, the convolutional operator applied in *strides*, skipping ahead by s steps after each feature. As we move up the hierarchy, each layer is s times smaller than the layer below it, effectively summarizing the input. This idea is shown in Figure 2.5. Multi-layer convolutional networks can also be augmented with “shortcut” connections, as in the ResNet model from § 2.2.2 (Johnson and Zhang, 2017).

Further reading

The deep learning textbook by Goodfellow et al. (2016) covers many of the topics in this chapter in more detail. For a comprehensive review of neural networks in natural language processing, see (Goldberg, 2017b). A seminal work on deep learning in natural language processing is the aggressively titled “Natural Language Processing (Almost) from Scratch”, which uses convolutional neural networks to perform a range of language processing tasks (Collobert et al., 2011b). This chapter focuses on feedforward and convolutional neural networks, but recurrent neural networks are one of the most important deep learning architectures for natural language processing. They are covered extensively in chapters 5 and 6.

The role of deep learning in natural language processing research has caused angst in some parts of the natural language processing research community (e.g., Goldberg, 2017a), especially as some of the more zealous deep learning advocates have argued that end-to-end learning from “raw” text can eliminate the need for linguistic constructs such as sentences, phrases, and even words (Zhang et al., 2015, originally titled *Text understanding from scratch*). These developments were surveyed by Manning (2016).

Exercises

1. Design a feedforward network to compute the XOR function:

$$f(x_1, x_2) = \begin{cases} -1, & x_1 = 1, x_2 = 1 \\ 1, & x_1 = 1, x_2 = 0 \\ 1, & x_1 = 0, x_2 = 1 \\ -1, & x_1 = 0, x_2 = 0 \end{cases} . \quad [2.57]$$

Your network should have a single output node which uses the Sign activation function. Use a single hidden layer, with ReLU activation functions. Describe all weights and offsets.

(c) Jacob Eisenstein 2014-2017. Work in progress.

2. Suppose that the parameters $\Theta = \{\Theta^{(x \rightarrow z)}, \Theta^{(z \rightarrow y)}, \mathbf{b}\}$ are a local optimum of a feedforward network in the following sense: there exists some $\epsilon > 0$ such that,

$$\begin{aligned} & \left(\|\tilde{\Theta}^{(x \rightarrow z)} - \Theta^{(x \rightarrow z)}\|_F^2 + \|\tilde{\Theta}^{(z \rightarrow y)} - \Theta^{(z \rightarrow y)}\|_F^2 + \|\tilde{\mathbf{b}} - \mathbf{b}\|_2^2 < \epsilon \right) \\ & \Rightarrow \left(L(\tilde{\Theta}) > L(\Theta) \right) \end{aligned} \quad [2.58]$$

Define the function π as a permutation on the hidden units, as described in § 2.3.3, so that for any Θ , $L(\Theta) = L(\Theta_\pi)$. Prove that if a feedforward network has a local optimum in the sense of Equation 2.58, then its loss is not a convex function of the parameters Θ , using the definition of convexity from § 1.3

3. Give the computational complexity of computing the convolutional operation described in Equation 2.52, in terms of all relevant dimensions.

Chapter 3

Linguistic applications of classification

Having learned some techniques for classification, we will now see how they can be applied to some classical problems in natural language technology. Later in this chapter, we discuss some of the design decisions involved in text classification, as well as evaluation practices.

3.1 Sentiment and opinion analysis

A popular application of text classification is to automatically determine 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 emotions are affected by phenomena such as the weather (Hannak et al., 2012), and how both opinions and 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).¹

Sentiment analysis can be framed as a fairly direct application of document classification, assuming 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

¹Comprehensive surveys on sentiment analysis and related problems are offered by Pang and Lee (2008) and Liu (2015).

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).

The bag-of-words model is a good fit for sentiment analysis at the document level: if the document is long enough, we would expect the words associated with its true sentiment to overwhelm the others. Indeed, **lexicon-based sentiment analysis** avoids machine learning altogether, and classifies documents by counting words against positive and negative sentiment word lists (Taboada et al., 2011).

The problem becomes more tricky for short documents, such as single-sentence reviews or social media posts. In these documents, linguistic issues like **negation** and **irrealis** (Polanyi and Zaenen, 2006) — events that are hypothetical or otherwise non-factual — can make bag-of-words classification ineffective. Consider the following examples:

- (3.1) That’s not bad for the first day.
- (3.2) This is not the worst thing that can happen.
- (3.3) It would be nice if you acted like you understood.
- (3.4) There is no reason at all to believe that the polluters are suddenly going to become reasonable. (Wilson et al., 2005)
- (3.5) This film should be brilliant. The actors are first grade. Stallone plays a happy, wonderful man. His sweet wife is beautiful and adores him. He has a fascinating gift for living life fully. It sounds like a great plot, **however**, the film is a failure. (Pang et al., 2002)

A minimal solution is to move from a bag-of-words model to a bag-of-**bigrams** model, where each base feature is a pair of adjacent words, e.g.,

$$\langle \text{that's}, \text{not} \rangle, \langle \text{not}, \text{bad} \rangle, \langle \text{bad}, \text{for} \rangle, \dots \quad [3.1]$$

Bigrams can handle relatively straightforward cases, such as when an adjective is immediately negated; trigrams would be required to extend to larger contexts (e.g., *not the worst*). But it should be clear that this approach will not scale to the more complex examples, such as (3.4) and (3.5). More sophisticated solutions try to account for the syntactic structure of the sentence (Wilson et al., 2005; Socher et al., 2013b), or apply more complex classifiers such as **convolutional neural networks** (Kim, 2014), which are described in chapter 2.

(c) Jacob Eisenstein 2014-2017. Work in progress.

3.1.1 Related problems

Subjectivity Closely related to sentiment analysis is **subjectivity detection**, which requires identifying the parts of a text that express subjective opinions, as well as other non-factual content such speculation and hypotheticals (Riloff and Wiebe, 2003). This can be done by treating each sentence as a separate document, and then applying a bag-of-words classifier: indeed, Pang and Lee (2004) do exactly this, using a training set consisting of (mostly) subjective sentences gathered from movie reviews, and (mostly) objective sentences gathered from plot descriptions. They augment this bag-of-words model with a graph-based algorithm that encourages nearby sentences to have the same subjectivity label.

Stance classification In debates, each participant takes a side: for example, advocating for or against adopting a vegetarian lifestyle or mandating free college education. The problem of stance classification involves identifying an author’s position from the text of the argument. In some cases, there is training data available for each position, so that standard document classification techniques can be employed. In other cases, it suffices to classify each document as whether it is in support or opposition of the argument advanced by a previous document (Anand et al., 2011). In the most challenging case, there is no labeled data for any of the stances, so the only possibility is group documents that advocate the same position (Somasundaran and Wiebe, 2009). This is a form of **unsupervised learning**, and will be discussed in chapter 4.

Targeted sentiment analysis The expression of sentiment is often more nuanced than a simple binary label. Consider the following examples:

(3.6) The vodka was good, but the meat was rotten.

(3.7) Go to Heaven for the climate, Hell for the company. (Mark Twain)

These statements display a mixed overall sentiment: positive towards some entities (e.g., *the vodka*), negative towards others (e.g., *the meat*). **Targeted sentiment analysis** seeks to identify the writer’s sentiment towards specific entities (Jiang et al., 2011). This requires identifying the entities in the text and linking them to specific sentiment words — much more than we can do with the classification-based approaches discussed thus far. For example, Kim and Hovy (2006) analyze sentence-internal structure to determine the topic of each sentiment expression.

Aspect-based opinion mining seeks to identify the sentiment of the author of a review towards predefined aspects such as PRICE and SERVICE, or, in the case of (3.7), CLIMATE and COMPANY (Hu and Liu, 2004). If the aspects are not defined in advance, it may again be necessary to employ **unsupervised learning** methods to identify them (e.g., Branavan et al., 2009).

(c) Jacob Eisenstein 2014-2017. Work in progress.

Emotion classification While sentiment analysis is framed in terms of positive and negative categories, psychologists generally regard **emotion** as more multifaceted. For example, Ekman (1992) argues that there are six basic emotions — happiness, surprise, fear, sadness, anger, and contempt — and that they are universal across human cultures. Alm et al. (2005) build a linear classifier for recognizing the emotions expressed in children’s stories. The ultimate goal of this work was to improve text-to-speech synthesis, so that stories could be read with intonation that reflected the emotional content. They used bag-of-words features, as well as features capturing the story type (e.g., jokes, folktales), and structural features that reflect the position of each sentence in the story. The task is difficult: even human annotators frequently disagreed with each other, and the best classifiers achieved accuracy between 60-70%.

3.1.2 Alternative approaches to sentiment analysis

Regression A more challenging version of sentiment analysis is to determine not just the class of a document, 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 $\hat{y} = \theta \cdot x + b$, where b is an offset. This approach is called **linear regression**, and sometimes **least squares**, because the regression coefficients θ are determined by minimizing the squared error, $(y - \hat{y})^2$. If the weights are regularized using a penalty $\lambda \|\theta\|_2^2$, then the name **ridge regression** is sometimes applied. Unlike logistic regression, both linear regression and ridge regression can be solved in closed form as a system of linear equations.

Ranking In many problems, the labels are ordered but discrete: for example, product reviews are often integers on a scale of 1 – 5, and grades are on a scale of $A - F$. Such **ranking** problems can be solved by discretizing the score $\theta \cdot x$ into ranks,

$$\hat{y} = \underset{r: \theta \cdot x \geq b_r}{\operatorname{argmin}} r, \quad [3.2]$$

where $b = [b_1 = -\infty, b_2, b_3, \dots, b_K]$ is a vector of boundaries. Crammer and Singer (2001) show that it is possible to learn both the weights and boundaries simultaneously, using a perceptron-like algorithm.

Lexicon-based classification Sentiment analysis is one of the only NLP tasks where hand-crafted feature weights are still widely employed. In **lexicon-based classification** (Taboada et al., 2011), the user creates a list of words for each label, and then classifies each document based on how many of the words from each list are present. In our linear classification framework, this is equivalent to choosing the following weights:

$$\theta_{y,j} = \begin{cases} 1, & j \in \mathcal{L}_y \\ 0, & \text{otherwise,} \end{cases} \quad [3.3]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

where \mathcal{L}_y is the lexicon for label y . Compared to the machine learning classifiers discussed in the previous chapters, lexicon-based classification may seem primitive. However, supervised machine learning relies on large annotated datasets, which are time-consuming and expensive to produce. If the goal is to distinguish two or more categories in a new domain, it may be simpler to start by writing down a list of words for each category.

An early lexicon was the *General Inquirer* (Stone, 1966). Today, popular sentiment lexicons include *sentiwordnet* (Esuli and Sebastiani, 2006) and an evolving set of lexicons from Liu (2015). For emotions and more fine-grained analysis, *Linguistic Inquiry and Word Count* (LIWC) provides a set of lexicons (Tausczik and Pennebaker, 2010). The MPQA lexicon indicates the polarity of some 8221 terms, as well as whether they are strongly or weakly subjective (Wiebe et al., 2005). A comprehensive comparison of sentiment lexicons is offered by Ribeiro et al. (2016). Given an initial **seed lexicon**, it is possible to automatically expand the lexicon by looking for words that frequently co-occur with words in the seed set (Hatzivassiloglou and McKeown, 1997; Qiu et al., 2011).

3.2 Word sense disambiguation

Consider the the following headlines:

(3.8) *Iraqi head seeks arms*

(3.9) *Prostitutes appeal to Pope*

(3.10) *Drunk gets nine years in violin case²*

These headlines 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. Word sense disambiguation is part of a larger field of research called **lexical semantics**, which is concerned with meanings of the words.

At a basic level, the problem of word sense disambiguation is to identify the correct sense for each word token in a document. Part-of-speech ambiguity (e.g., noun versus verb) is usually considered to be a different problem, to be solved at an earlier stage. From a linguistic perspective, senses are not really properties of words, but of **lemmas**, which are canonical forms that stand in for a set of inflected words. For example, *arm*/N is a lemma that includes the inflected form *arms*/N — the /N indicates that it we are referring to the noun form of the word. Similarly, *arm*/V is a lemma that includes the inflected verbs (*arm*/V, *arms*/V, *armed*/V, *arming*/V). Therefore, word sense disambiguation requires first identifying the correct part-of-speech and lemma for each token, and

²These examples, and many more, can be found at <http://www.ling.upenn.edu/~beatrice/humor/headlines.html>

then choosing the correct sense from the inventory associated with the corresponding lemma.³

3.2.1 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*
- [CONTRIBUTE TO]: *His evasive replies only served to heighten suspicion*
- [PROVIDE]: *We serve only the rawest fish*
- [ENLIST]: *She served in an elite combat unit*
- [JAIL]: *He served six years for a crime he didn't commit*
- [LEGAL]: *They were served with subpoenas*⁴

These sense distinctions are annotated in WORDNET (<http://wordnet.princeton.edu>), a lexical semantic database for English. WordNet consists of roughly 100,000 **synsets**, which are groups of lemmas (or phrases) that are synonymous. An example synset is $\{chump^1, fool^2, sucker^1, mark^9\}$, where the superscripts index the sense of each lemma that is included in the synset: for example, there are at least eight other senses of *mark* that have different meanings, and are not part of this synset. A lemma is **polysemous** if it participates in multiple synsets.

WordNet plays an key role in setting the parameters of the word sense disambiguation problem, and in formalizing lexical semantic knowledge of English. (WordNets have been created for a few dozen other languages, at varying levels of detail.) 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 (Kilgariff, 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), people agree on roughly 70% of examples using WordNet senses; far better than chance, but perhaps less than we might like.

***Other lexical semantic relations** Besides **synonymy**, WordNet also describes many other lexical semantic relationships, including:

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

³Navigli (2009) provides a survey of approaches for word-sense disambiguation.

⁴Several of the examples are adapted from WordNet (Fellbaum, 2010)

- **hyponymy**: x is a special case of y , e.g. RED-COLOR; the inverse relationship is **hypernymy**;
- **meronymy**: x is a part of y , e.g., WHEEL-BICYCLE; the inverse relationship is **holonymy**.

Classification of these relations can be performed by searching for characteristic patterns between pairs of words, e.g., X , *such as* Y , which signals hyponymy (Hearst, 1992), or X *but* Y , which signals antonymy (Hatzivassiloglou and McKeown, 1997). Another approach is to analyze each term's **distributional statistics** (the frequency of its neighboring words). Such approaches are described in detail in chapter 13.

3.2.2 Word sense disambiguation as classification

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

- (3.11) Town officials are hoping to attract new manufacturing plants through weakened environmental regulations.
- (3.12) The endangered plants play an important role in the local ecosystem.

It is possible to build a feature vector using the bag-of-words representation, by treating each context as a pseudo-document. We can then construct a feature function for each potential sense y ,

$$f(\langle \text{plant}, \text{The endangered plants play an } \dots \rangle, y) = \{ \langle \text{the}, y \rangle : 1, \langle \text{endangered}, y \rangle : 1, \langle \text{play}, y \rangle : 1, \langle \text{an}, y \rangle : 1, \dots \}$$

As in document classification, many of these features are irrelevant, but a few are very strong indicators. In this example, the context word *endangered* is a strong signal that the intended sense is biology rather than manufacturing. We would therefore expect a learning algorithm to assign high weight to $\langle \text{endangered}, \text{BIOLOGY} \rangle$, and low weight to $\langle \text{endangered}, \text{MANUFACTURING} \rangle$.⁵

It may also be helpful to go beyond the bag-of-words: for example, one might encode the position of each context word with respect to the target, e.g.,

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

⁵The context bag-of-words can be also used to perform word-sense disambiguation without machine learning: the Lesk (1986) algorithm selects the word sense whose dictionary definition best overlaps the local context.

These **collocation features** give more information about the specific role played by each context word. This idea can be taken further by incorporating additional syntactic information about the grammatical role played by each context feature, such as the **dependency path** (see chapter 10).

After deciding on the features, we can train a classifier to predict the sense of each word. 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, annotated as part of the WordNet project (Fellbaum, 2010). SemCor annotations look like this:

(3.13) As of Sunday_{*n*}¹ night_{*n*}¹ there was_{*v*}⁴ no word_{*n*}² . . . ,

with the superscripts indicating the annotated sense of each polysemous word.

As always, supervised classification is only possible if enough labeled examples can be accumulated. This is difficult in word sense disambiguation, because each polysemous lemma requires its own training set: having a good classifier for the senses of *serve* is no help towards disambiguating *plant*. For this reason, **unsupervised** and **semisupervised** methods are particularly important for WSD (e.g., Yarowsky, 1995). These methods will be discussed in chapter 4. Unsupervised methods typically lean on the heuristic of “one sense per discourse”, which means that a lemma will usually have a single, consistent sense throughout any given document (Gale et al., 1992). Based on this heuristic, we can propagate information from high-confidence instances to lower-confidence instances in the same document (Yarowsky, 1995).

3.3 Design decisions for text classification

Text classification involves a number of design decisions. Some of these decisions, such as smoothing or regularization, are specific to the classifier. These decisions are described in the previous chapter. But even the construction of the feature vector itself involves a number of design decisions, and these decisions can be the most consequential for the classifier’s performance.

3.3.1 What is a word?

The bag-of-words representation presupposes that extracting a vector of word counts from text is unambiguous. But text documents are generally represented as a sequences of characters, and the conversion to bag-of-words presupposes a definition of the “words” that are to be counted.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Whitespace	Isn't	Ahab,	Ahab?	;))					
Treebank	Is	n't	Ahab	,	Ahab	?	;)	
Tweet	Isn't	Ahab	,	Ahab	?	;))			
TokTok (Dehdari, 2014)	Isn	'	t	Ahab	,	Ahab	?	;)

Figure 3.1: The output of four `nltk` tokenizers, applied to the string *Isn't Ahab, Ahab? ;)*

Tokenization

The first subtask for constructing a bag-of-words vector is **tokenization**: converting the text from a sequence of characters to a sequence of **word tokens**. A simple approach is to define a subset of characters as whitespace, and then split the text on these tokens. However, whitespace-based tokenization is not ideal: we may want to split conjunctions like *isn't* and hyphenated phrases like *prize-winning* and *half-asleep*, and we likely want to separate words from commas and periods that immediately follow them. This suggests that a tokenizer should split on all non-alphanumeric characters, but we would prefer not to split abbreviations like *U.S.* and *Ph.D.* In languages with Roman scripts, tokenization is typically performed using regular expressions, with modules designed to handle each of these cases. For example, the `nltk` package includes a number of tokenizers; the outputs of four of the better-known tokenizers are shown in Figure 3.1. Social media researchers have found that emoticons and other forms of orthographic variation pose new challenges for tokenization, leading to the development of special purpose tokenizers to handle these phenomena (O'Connor et al., 2010).

Tokenization is a language-specific problem, and each language poses unique challenges. For example, Chinese does not include spaces between words, nor any other consistent orthographic markers of word boundaries. A “greedy” approach is to scan the input for character substrings that are in a predefined lexicon. However, Xue et al. (2003) notes that this can be ambiguous, since many character sequences could be segmented in multiple ways. Instead, he trains a classifier to determine whether each Chinese character, or **hanzi**, is a word boundary. More advanced sequence labeling methods for word segmentation are discussed in § 7.4. Similar problems can occur in languages with alphabetic scripts, such as German, which does not include whitespace in compound nouns, yielding examples such as *Freundschaftsbezeugungen* and *Dilettantenaufdringlichkeiten* [**todo: ask German speaker for better examples**]. As Twain (1997) argues, “*These things are not words, they are alphabetic processions.*” Social media raises similar problems for English and other languages, with hashtags such as *#TrueLoveInFourWords* requiring decomposition for analysis (Brun and Roux, 2014).

(c) Jacob Eisenstein 2014-2017. Work in progress.

Original	The	Williams	sisters	are	leaving	this	tennis	centre
Porter stemmer	the	william	sister	are	leav	thi	tenni	centr
Lancaster stemmer	the	william	sist	ar	leav	thi	ten	cent

Figure 3.2: The outputs of the Porter (1980) and Lancaster (Paice, 1990) stemmers

Normalization

After splitting the text into tokens, the next question is which tokens are really distinct. Is it necessary to distinguish *great*, *Great*, and *GREAT*? Sentence-initial capitalization may be irrelevant to the classification task. The elimination of case distinctions will result in a smaller vocabulary, and thus smaller feature vectors. However, case distinctions might be relevant in some situations: for example, *apple* is a delicious pie filling, while *Apple* is a company dedicated to marketing proprietary dongles and power adapters. For Latin script, case conversion can be performed using unicode string libraries. Many scripts do not have case distinctions (e.g., the Devanagari script used for South Asian languages, the Thai alphabet, and Japanese kana), and case conversion for all scripts may not be available in every programming environment.

Case conversion is an example of **normalization**, which refers to string transformations that remove distinctions that are felt to be irrelevant (Sproat et al., 2001). Other normalizations include the standardization of numbers (e.g., *1,000* to *1000*) and dates (e.g., *August 11, 2015* to *2015/11/08*). Depending on the application, it may even be worthwhile to convert all numbers and dates to special tokens, `!NUM` and `!DATE`. Social media features orthographic phenomena such as expressive lengthening (e.g., *coooooool*), which may also be normalized (Aw et al., 2006; Yang and Eisenstein, 2013). Similarly, historical texts feature spelling variations which may be normalized to a standard form (Baron and Rayson, 2008).

A more extreme form of normalization is to eliminate **inflectional affixes**, such as the *-ed* and *-s* suffixes in English. On this view, *bike*, *bikes*, *biking*, and *biked* all refer to the same underlying concept, so they should be grouped into a single feature. A **stemmer** is a program for eliminating affixes, usually by applying a series of regular expression substitutions. Character-based stemming algorithms are necessarily approximate, as shown in Figure 3.2: the Lancaster stemmer incorrectly identifies *-ers* as an inflectional suffix of *sisters* (by analogy to *fix/fixers*), and both stemmers incorrectly identify *-s* as a suffix of *this* and *Williams*. Fortunately, even inaccurate stemming can improve bag-of-words classification models, by merging related strings and thereby reducing the vocabulary size.

Accurately handling irregular orthography requires word-specific rules. **Lemmatizers** are systems that identify the underlying lemma of a given wordform. They must avoid the over-generalization errors of the stemmers in Figure 3.2, and also handle more complex



Figure 3.3: Tradeoff between token coverage (y-axis) and vocabulary size, on the `nltk` movie review dataset, after sorting the vocabulary by decreasing frequency. The red dashed lines indicate 80%, 90%, and 95% coverage.

transformations, such as *geese* \rightarrow *goose*. The output of the WordNet lemmatizer is shown in the final line of Figure 3.2. Both stemming and lemmatization are language-specific: an English stemmer or lemmatizer is of little use on a text written in another language. The discipline of **morphology** relates to the study of word-internal structure, and is described in more detail in § 8.1.2.

The value of normalization depends on the data and the task. Normalization reduces the size of the feature space, which can help in generalization. However, there is always the risk of merging away linguistically meaningful distinctions. In supervised machine learning, regularization and smoothing can play a similar role to normalization — preventing the learner from overfitting to rare features — while avoiding the language-specific engineering required for accurate normalization. In unsupervised scenarios, such as content-based information retrieval (Manning et al., 2008) and topic modeling (Blei et al., 2003), normalization is more critical.

3.3.2 How many words?

Limiting the size of the feature vector reduces the memory footprint of the resulting models, and increases the speed of prediction. Normalization can help to play this role, but a more direct approach is simply to limit the vocabulary to the N most frequent words in the dataset. For example, in the `movie-reviews` dataset provided with `nltk` (originally from Pang et al., 2002), there are 39,768 word types, and 1.58M tokens. As shown in Figure 3.3a, the most frequent 4000 word types cover 90% of all tokens, offering an order-of-magnitude reduction in the model size. Such ratios are language-specific: in for example, in the Brazilian Portuguese Mac-Morpho corpus (Aluísio et al., 2003), attaining 90% coverage requires more than 10000 word types (Figure 3.3b). This reflects the morphological complexity of Portuguese, which includes many more inflectional suffixes than English.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Eliminating rare words is not always advantageous for classification performance: for example, names, which are typically rare, play a large role in distinguishing topics of news articles. Another way to reduce the size of the feature space is to eliminate **stopwords** such as *the*, *to*, and *and*, which may seem to play little role in expressing the topic, sentiment, or stance. This is typically done by creating a **stoplist** (e.g., `nltk.corpus.stopwords`), and then ignoring all terms that match the list. However, corpus linguists and social psychologists have shown that seemingly inconsequential words can offer surprising insights about the author or nature of the text (Biber, 1991; Chung and Pennebaker, 2007). Furthermore, high-frequency words are unlikely to cause overfitting in well-regularized discriminative classifiers. As with normalization, stopwords filtering is more important for unsupervised problems, such as term-based document retrieval; in this case, matching a *to* or *the* in the search query offers little information that the resulting document will meet the goals of the user’s search.

Another alternative for controlling model size is **feature hashing** (Weinberger et al., 2009). Each feature is assigned an index using a hash function. If a hash function that permits collisions is chosen (typically by taking the hash output modulo some integer), then the model can be made arbitrarily small, as multiple features share a single weight. Because most features are rare, accuracy is surprisingly robust to such collisions (Ganchev and Dredze, 2008).

3.3.3 Count or binary?

Finally, we may consider whether we want our feature vector to include the **count** of each word, or its **presence**. This gets at a subtle limitation of linear classification: two *failures* may be worse than one, but is it really twice as bad? Motivated by this intuition, Pang et al. (2002) use binary indicators of presence or absence in the feature vector: $f_j(\mathbf{x}, y) \in \{0, 1\}$. They find that classifiers trained on these binary vectors tend to outperform feature vectors based on word counts. One explanation is that words tend to appear in clumps: if a word has appeared once in a document, it is likely to appear again (Church, 2000). These subsequent appearances can be attributed to this tendency towards repetition, and thus provide little additional information about the class label of the document.

3.4 Evaluating classifiers

In any supervised machine learning application, it is critical to reserve a held-out test set, and use this data for only one purpose: to evaluate the overall accuracy of a single classifier. Using this data more than once would cause the estimated accuracy to be overly optimistic, because the classifier would be customized to this data, and would not perform as well as on unseen data in the future. It is usually necessary to set hyperparameters or

perform feature selection, so you may need to construct a **tuning** or **development set** for this purpose, as discussed in § 1.1.5.

There are a number of ways to evaluate classifier performance. The simplest is **accuracy**: the number of correct predictions, divided by the total number of instances,

$$\text{acc}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{N} \sum_i^N \delta(y^{(i)} = \hat{y}). \quad [3.4]$$

If you've ever taken an exam, it was probably graded by accuracy. Why are other metrics necessary? The main reason is **class imbalance**. Suppose we were building a classifier to detect whether a electronic health record (EHR) described symptoms of a rare disease, which appear in only 1% of all documents in the dataset. A classifier that reports $\hat{y} = -1$ for all documents would achieve 99% accuracy, but would be practically useless. We need metrics that are capable of detecting the classifier's ability to discriminate between classes, even when the distribution is skewed.

One solution would be to build a **balanced test set**, in which 50% of documents are positive. But this would mean throwing away 98% of the original dataset! Furthermore, the detection threshold itself might be a design consideration: in health-related applications, we might prefer a very sensitive classifier, which returned a positive prediction if there is even a small chance that $y^{(i)} = +1$. In other applications, a positive result might trigger a costly action, so we would prefer a classifier that only makes positive predictions when absolutely certain. We need additional metrics to capture these characteristics.

3.4.1 Precision, recall, and F -measure

For any label (e.g., positive for presence of symptoms of a disease), there are two possible errors:

- **False positive**: the system incorrectly predicts the label.
- **False negative**: the system incorrectly fails to predict the label.

Similarly, for any label, there are two ways to be correct:

- **True positive**: the system correctly predicts the label.
- **True negative**: the system correctly predicts that the label does not apply to this instance.

Classifiers that make a lot of false positive errors are too sensitive; classifiers that make a lot of false negative errors are not sensitive enough. We can capture these two behaviors

through two additional metrics, **recall** and **precision**:

$$\text{recall}(\mathbf{y}, \hat{\mathbf{y}}, k) = \frac{\text{TP}}{\text{TP} + \text{FN}} \quad [3.5]$$

$$\text{precision}(\mathbf{y}, \hat{\mathbf{y}}, k) = \frac{\text{TP}}{\text{TP} + \text{FP}}. \quad [3.6]$$

Recall and precision are both conditional likelihoods of a correct prediction, which is why their numerators are the same. Recall is conditioned on k being the correct label, $y^{(i)} = k$, so the denominator sums over true positive and false negatives. Precision is conditioned on k being the prediction, so the denominator sums over true positives and false positives. Note that true negatives are not considered in either statistic. The classifier that labels every document as “negative” would achieve zero recall; precision would be $\frac{0}{0}$.

Recall and precision are complementary objectives. A high-recall classifier is preferred when false negatives are cheaper than false positives: for example, in a preliminary screening for symptoms of a disease, the cost of a false positive might be an additional test, while a false negative would result in the disease going untreated. Conversely, we prefer a high-precision classifier when false positives are more expensive: for example, in spam detection, a false negative is a relatively minor inconvenience, while a false positive might mean that an important message goes unread.

The **F -measure** combines recall and precision into a single metric, using the harmonic mean:

$$F\text{-measure}(\mathbf{y}, \hat{\mathbf{y}}, k) = \frac{2rp}{r + p}, \quad [3.7]$$

where r is recall and p is precision.⁶

Evaluating multi-class classification Recall, precision, and F -measure are defined with respect to a specific label k . When there are multiple labels of interest (e.g., in word sense disambiguation), it is necessary to combine the F -measure across each class. **Macro F -measure** is the average F -measure across several classes,

$$\text{Macro-}F(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{|\mathcal{K}|} \sum_{k \in \mathcal{K}} F\text{-measure}(\mathbf{y}, \hat{\mathbf{y}}, k) \quad [3.8]$$

In multi-class problems with unbalanced class distributions, the macro F -measure is a balanced measure of how well the classifier recognizes each class. In **micro F -measure**, we compute true positives, false positives, and false negatives for each class, and then add them up to compute a single recall, precision, and F -measure. This metric is balanced across instances rather than classes, so it weights each class in proportion to its frequency — unlike macro F -measure, which weights each class equally.

⁶ F -measure is sometimes called F_1 , and generalizes to $F_\beta = \frac{(1+\beta^2)rp}{\beta^2 p + r}$. The β parameter can be tuned to emphasize recall or precision.

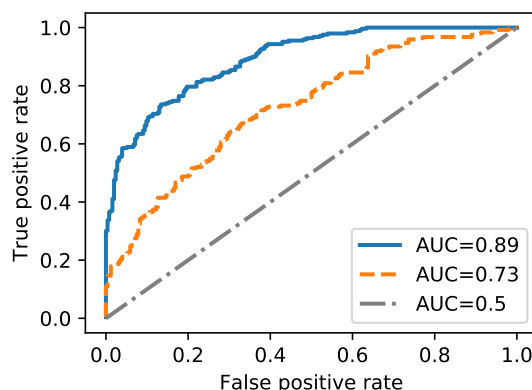


Figure 3.4: ROC curves for three classifiers of varying discriminative power, measured by AUC (area under the curve)

3.4.2 Threshold-free metrics

In binary classification problems, it is possible to trade off between recall and precision by adding a constant “threshold” to the output of the scoring function. This makes it possible to trace out a curve, where each point indicates the performance at a single threshold. In the **receiver operating characteristic (ROC) curve**,⁷ the x-axis indicates the **false positive rate**, $\frac{FP}{FP+TN}$, and the y-axis indicates the recall, or **true positive rate**. A perfect classifier attains perfect recall without any false positives, tracing a “curve” from the origin (0,0) to the upper left corner (0,1), and then to (1,1). In expectation, a non-discriminative classifier traces a diagonal line from the origin (0,0) to the upper right corner (1,1). Real classifiers tend to fall between these two extremes. Examples are shown in Figure 3.4.

The ROC curve can be summarized in a single number by taking its integral, the **area under the curve (AUC)**. The AUC has an intuitive interpretation as the probability that a randomly-selected positive example will be assigned a higher score by the classifier than a randomly-selected negative example. Thus, a perfect classifier has $AUC = 1$ (all positive examples score higher than all negative examples); a random non-discriminative classifier has $AUC = 0.5$ (given a randomly selected positive and negative example, either could score higher with equal probability); a perfectly wrong classifier would have $AUC = 0$ (all negative examples score higher than all positive examples). One advantage of AUC in comparison to F -measure is that the baseline rate of 0.5 does not depend on the label distribution.

⁷The name “receiver operator characteristic” comes from the metric’s origin in signal processing applications (Peterson et al., 1954). Other threshold-free metrics include precision-recall curves, precision-at- k , and balanced F -measure; see Manning et al. (2008) for more details.

3.4.3 Classifier comparison and statistical significance

Building NLP systems often involves comparing different classification techniques. In some cases, the comparison is between algorithms, such as logistic regression versus averaged perceptron, or L_2 regularization versus L_1 . In other cases, the comparison is between feature sets, such as the bag-of-words versus positional bag-of-words feature sets discussed in § 3.2.2. **Ablation testing** involves systematically removing (ablating) various aspects of the classifier, such as feature groups, and testing the **null hypothesis** that the ablated classifier is as good as the full model.

A full treatment of hypothesis testing is beyond the scope of this text, but this section contains a brief summary of the techniques necessary to compare classifiers. The main aim of hypothesis testing is to determine whether the difference between two statistics — for example, the accuracies of two classifiers — is likely to arise by chance. We will be concerned with chance fluctuations that arise due to the finite size of the test set.⁸ An improvement of 10% on a test set with ten instances may reflect a random fluctuation that makes the test set more favorable to classifier c_1 than c_2 ; on another test set with a different ten instances, we might find that c_2 does better than c_1 . But if we observe the same 10% improvement on a test set with 1000 instances, this is highly unlikely to be explained by chance. Such a finding is said to be **statistically significant** at a level p , which is the probability of an effect of equal or greater magnitude when the null hypothesis (that the classifiers are equally accurate) is true. For example, we write $p < .05$ when the likelihood of an equal or greater effect is less than 5%, assuming the null hypothesis is true.⁹

The binomial test

The statistical significance of a difference in accuracy can be evaluated using classical tests, such as the **binomial test**.¹⁰ Suppose that classifiers c_1 and c_2 disagree on N instances in the test set, and that c_1 is correct on k of those instances. Under the null hypothesis that the classifiers are equally accurate, we would expect k/N to be roughly equal to $1/2$. As N increases, k/N should be increasingly close to $1/2$. These properties are captured by the **binomial distribution**, which is a probability over counts of binary random variables.

⁸Other sources of variance include the initialization of non-convex classifiers such as neural networks, and the ordering of instances in online learning such as stochastic gradient descent and perceptron.

⁹Statistical hypothesis testing is useful only to the extent that the existing test set is representative of the instances that will be encountered in the future. If, for example, the test set is constructed from news documents, no hypothesis test can predict which classifier will perform best on documents from another domain, such as electronic health records.

¹⁰A well-known alternative to the binomial test is **McNemar's test**, which computes a **test statistic** based on the number of examples that are correctly classified by one system and incorrectly classified by the other. The null hypothesis distribution for this test statistic is known to be drawn from a chi-squared distribution with a single degree of freedom, so a p -value can be computed from the cumulative density function of this distribution (Dietterich, 1998). Both tests give similar results in most circumstances, but the binomial test is easier to explain from first principles.

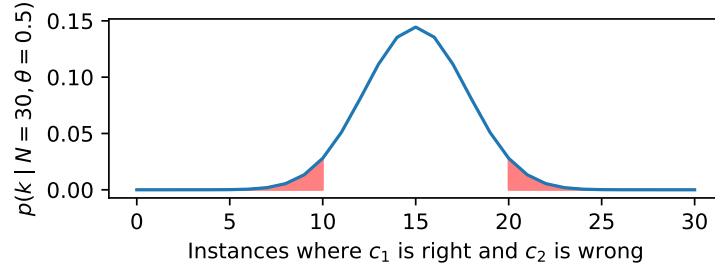


Figure 3.5: Probability mass function for the binomial distribution. The pink highlighted areas represent the cumulative probability for a significance test on an observation of $k = 10$ and $N = 30$.

We write $k \sim \text{Binom}(\theta, N)$ to indicate that k is drawn from a binomial distribution, with parameter N indicating the number of random “draws”, and θ indicating the probability of “success” on each draw. The **probability mass function** (PMF) of the binomial distribution is,

$$p_{\text{Binom}}(k; N, \theta) = \binom{N}{k} \theta^k (1 - \theta)^{N-k}, \quad [3.9]$$

with θ^k representing the probability of the k successes, $(1 - \theta)^{N-k}$ representing the probability of the $N - k$ unsuccessful draws. The expression $\binom{N}{k} = \frac{N!}{k!(N-k)!}$ is a binomial coefficient, representing the number of possible orderings of events; this ensures that the distribution sums to one over all $k \in \{0, 1, 2, \dots, N\}$.

Under the null hypothesis, $\theta = \frac{1}{2}$: when the classifiers disagree, they are each equally likely to be right. Now suppose that among N disagreements, c_1 is correct only $k < \frac{N}{2}$ times. The probability of c_1 being correct k or fewer times is the **one-tailed p-value**, because it is computed from the area under the binomial probability mass function from 0 to k , as shown in the left tail of Figure 3.5. This **cumulative probability** is computed as a sum over all values $i \leq k$,

$$\Pr_{\text{Binom}} \left(\text{count}(\hat{y}_2^{(i)} = y^{(i)} \neq \hat{y}_1^{(i)}) \leq k; N, \theta = \frac{1}{2} \right) = \sum_{i=0}^k p_{\text{Binom}} \left(i; N, \theta = \frac{1}{2} \right). \quad [3.10]$$

The one-tailed p-value applies only to the asymmetric null hypothesis that c_1 is at least as accurate as c_2 . To test the **two-tailed** null hypothesis that c_1 and c_2 are equally accurate, we would take the sum of one-tailed p-values, where the second term is computed from the right tail of Figure 3.5. The binomial distribution is symmetric, so this can be computed by simply doubling the one-tailed p-value.

Algorithm 7 Bootstrap sampling for classifier evaluation. The original test set is $\{\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)}\}$, the metric is $\delta(\cdot)$, and the number of samples is M .

```

procedure BOOTSTRAP-SAMPLE( $\mathbf{x}^{(1:N)}, \mathbf{y}^{(1:N)}, \delta(\cdot), M$ )
  for  $t \in \{1, 2, \dots, M\}$  do
    for  $i \in \{1, 2, \dots, N\}$  do
       $j \sim \text{UniformInteger}(1, N)$ 
       $\tilde{\mathbf{x}}^{(i)} \leftarrow \mathbf{x}^{(j)}$ 
       $\tilde{\mathbf{y}}^{(i)} \leftarrow \mathbf{y}^{(j)}$ 
     $d^{(t)} \leftarrow \delta(\tilde{\mathbf{x}}^{(1:N)}, \tilde{\mathbf{y}}^{(1:N)})$ 
  return  $\{d^{(t)}\}_{t=1}^M$ 

```

Two-tailed tests are more stringent, but they are necessary in cases in which there is no prior intuition about whether c_1 or c_2 is better. For example, in comparing logistic regression versus averaged perceptron, a two-tailed test is appropriate. In an ablation test, c_2 may contain a superset of the features available to c_1 . If the additional features are thought to be likely to improve performance, then a one-tailed test would be appropriate, if chosen in advance. However, such a test can only prove that c_2 is more accurate than c_1 , and not the reverse.

*Randomized testing

The binomial test is appropriate for accuracy, but not for more complex metrics such as F -measure. To compute statistical significance for arbitrary metrics, we must turn to randomization. Specifically, we draw a set of M **bootstrap samples** (Efron and Tibshirani, 1993), by resampling instances from the original test set with replacement. Each bootstrap sample is itself a test set of size N . Some instances from the original test set will not appear in any given bootstrap sample, while others will appear multiple times; but overall, the sample will be drawn from the same distribution as the original test set. We can then compute any desired evaluation on each bootstrap sample, which gives a distribution over the value of the metric. Algorithm 7 shows how to perform this computation.

To compare the F -measure of two classifiers c_1 and c_2 , we set the function $\delta(\cdot)$ to compute the difference in F -measure on the bootstrap sample. If the difference is less than or equal to zero in at least 5% of the samples, then we cannot reject the one-tailed null hypothesis that c_2 is at least as good as c_1 (Berg-Kirkpatrick et al., 2012). We may also be interested in the 95% **confidence interval** around a metric of interest, such as the F -measure of a single classifier. This can be computed by sorting the output of Algorithm 7, and then setting the top and bottom of the 95% confidence interval to the values at the 2.5% and 97.5% percentiles of the sorted outputs. Alternatively, you can fit a normal distribution to the set of differences across bootstrap samples, and compute a Gaussian

(c) Jacob Eisenstein 2014-2017. Work in progress.

confidence interval from the mean and variance.

As the number of bootstrap samples goes to infinity, $M \rightarrow \infty$, the bootstrap estimate is increasingly accurate. A typical choice for M is 10^4 or 10^5 ; larger values are necessary for smaller p -values. One way to validate your choice of M is to run the test multiple times, and ensure that the p -values are similar; if not, increase M by an order of magnitude. This is a heuristic measure of the **variance** of the test, which can decrease with the square root \sqrt{M} (Robert and Casella, 2013).

3.4.4 *Multiple comparisons

Sometimes it is necessary to perform multiple hypothesis tests, such as when comparing the performance of several classifiers on multiple datasets. Suppose you have five datasets, and you compare four versions of your classifier against a baseline system, for a total of 20 comparisons. Even if none of your classifiers is better than the baseline, there will be some chance variation in the results, and in expectation you will get one statistically significant improvement at $p = 0.05 = \frac{1}{20}$. It is therefore necessary to adjust the p -values when reporting the results of multiple comparisons.

One approach is to require a threshold of $\frac{\alpha}{m}$ to report a p value of $p < \alpha$ when performing m tests. This is known as the **Bonferroni correction**, and it limits the overall probability of incorrectly rejecting the null hypothesis at α . Another approach is to bound the **false discovery rate** (FDR), which is the fraction of null hypothesis rejections that are incorrect. Benjamini and Hochberg (1995) propose a p -value correction that bounds the fraction of false discoveries at α : sort the p -values of each individual test in ascending order, and set the significance threshold equal to largest k such that $p_k \leq \frac{k}{m}\alpha$. If $k > 1$, the FDR adjustment is more permissive than the Bonferroni correction.

3.5 Building datasets

For many classification tasks of interest, no labeled data exists. If you want to build a classifier, you must first build a dataset of your own. This includes selecting a set of documents or instances to annotate, and then performing the annotations.

In many cases, the scope of the dataset is determined by the application: if you want to build a system to classify electronic health records, then you must work with a corpus of records of the type that your classifier will encounter when deployed. In other cases, the goal is to build a system that will work across a broad range of documents. In this case, it is best to have a **balanced** corpus, with contributions from many styles and genres. For example, the Brown corpus draws from texts ranging from government documents to romance novels (Francis, 1964), and the Google Web Treebank includes annotations for five “domains” of web documents: question answers, emails, newsgroups, reviews, and blogs (Petrov and McDonald, 2012).

(c) Jacob Eisenstein 2014-2017. Work in progress.

3.5.1 Metadata as labels

Annotation is difficult and time-consuming, and most people would rather avoid it. Luckily, it is sometimes possible to exploit existing metadata to obtain the desired labels. For example, reviews are often accompanied by a numerical rating, which can be converted into a classification label (see § 3.1). Similarly, the nationalities of social media users can be estimated from their profiles (Dredze et al., 2013) or even the time zones of their posts (Gouws et al., 2011). More ambitiously, we may try to classify the political affiliations of social media profiles based on their social network connections to politicians and major political parties (Rao et al., 2010).

The convenience of quickly constructing large labeled datasets without manual annotation is appealing. However this approach relies on the assumption that unlabeled instances — for which metadata is unavailable — will be similar to labeled instances. Consider the example of labeling the political affiliation of social media users based on their network ties to politicians. If a classifier attains high accuracy on such a test set, can we assume that it accurately predicts the political affiliation of all social media users? Probably not. Social media users who establish social network ties to politicians may also be more likely to mention politics in the text of their messages, as compared to the average user, for whom no political metadata is available. If so, the accuracy on a test set constructed from social network metadata would give an overly optimistic picture of the method's true performance on unlabeled data.

3.5.2 Labeling data

In many cases, there is no way to get ground truth labels other than manual annotation. Good annotations should satisfy several criteria: they should be **expressive** enough to capture the phenomenon of interest; they should be **replicable**, meaning that another annotator or team of annotators would produce very similar annotations if given the same data; and they should be **scalable**, so that they can be produced relatively quickly. Hovy and Lavid (2010) propose a structured procedure for obtaining annotations that meet these criteria, which is summarized below.

1. **Determine what the annotations are to include.** This is usually based on some theory of the underlying phenomenon: for example, if the goal is to produce annotations about the emotional state of a document's author, one should start with an theoretical account of the types or dimensions of emotion (e.g., Mohammad and Turney, 2013). At this stage, the tradeoff between expressiveness and scalability should be considered: a full instantiation of the underlying theory might be too costly to annotate at scale, so reasonable approximations should be considered.
2. Optionally, one may **design or select a software tool to support the annotation**

(c) Jacob Eisenstein 2014-2017. Work in progress.

effort. Existing general-purpose annotation tools include BRAT (Stenetorp et al., 2012) and MMAX2 (Müller and Strube, 2006).

3. **Formalize the instructions for the annotation task.** To the extent that the instructions are not explicit, the resulting annotations will depend on the intuitions of the annotators. These intuitions may not be shared by other annotators, or by the users of the annotated data. Therefore explicit instructions are critical to ensuring the annotations are replicable and usable by other researchers.
4. **Perform a pilot annotation** of a small subset of data, with multiple annotators for each instance. This will give a preliminary assessment of both the replicability and scalability of the current annotation instructions. Metrics for computing the rate of agreement are described below. Manual analysis of specific disagreements should help to clarify the instructions, and may lead to modifications of the annotation task itself. For example, if two labels are commonly conflated by annotators, it may be best just to merge them.
5. After finalizing the annotation protocol and instructions, the main annotation effort can begin. Some if not all of the instances should receive multiple annotations, so that inter-annotator agreement can be computed. In some annotation projects, instances receive many annotations, which are then aggregated into a “consensus” label (e.g., Danescu-Niculescu-Mizil et al., 2013). However, if the annotations are time-consuming or require significant expertise, it may be preferable to maximize scalability by obtaining multiple annotations for only a small subset of examples.
6. Compute and report inter-annotator agreement, and release the data. In some cases, the raw text data cannot be released, usually due to concerns related to copyright or privacy. In these cases, one solution is to publicly release **stand-off annotations**, which contain links to document identifiers. The documents themselves can be released under the terms of a licensing agreement.

Measuring inter-annotator agreement

To measure the replicability of annotations, a standard practice is to compute the extent to which annotators agree with each other. If the annotators frequently disagree, this casts doubt on either their reliability or on the annotation system itself. For classification, one can compute the frequency with which the annotators agree; for rating scales, one can compute the average distance between ratings. These raw agreement statistics must then be compared with the rate of **chance agreement** — the level of agreement that would be obtained between two annotators who ignored the data.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Cohen’s Kappa is widely used for quantifying the agreement on discrete labeling tasks (Cohen, 1960; Carletta, 1996),¹¹

$$\kappa = \frac{\text{agreement} - E[\text{agreement}]}{1 - E[\text{agreement}]}.$$
 [3.11]

The numerator is the difference between the observed agreement and the chance agreement, and the denominator is the difference between perfect agreement and chance agreement. Thus, $\kappa = 1$ when the annotators agree in every case, and $\kappa = 0$ when the annotators agree only as often as would happen by chance. Various heuristic scales have been proposed for determining when κ indicates “moderate”, “good”, or “substantial” agreement; for reference, Lee and Narayanan (2005) report $\kappa \approx 0.45 - 0.47$ for annotations of emotions in spoken dialogues, which they describe as “moderate agreement”; Stolcke et al. (2000) report $\kappa = 0.8$ for annotations of **dialogue acts**, which are labels for the purpose of each turn in a conversation.

When there are two annotators, the expected chance agreement is computed as,

$$E[\text{agreement}] = \sum_k \hat{\text{Pr}}(Y = k)^2,$$
 [3.12]

where k is a sum over labels, and $\hat{\text{Pr}}(Y = k)$ is the empirical probability of label k across all annotations. The formula is derived from the expected number of agreements if the annotations were randomly shuffled. Thus, in a binary labeling task, if one label is applied to 90% of instances, chance agreement is $.9^2 + .1^2 = .82$.

Crowdsourcing

Crowdsourcing is often used to rapidly obtain annotations for classification problems. For example, **Amazon Mechanical Turk** makes it possible to define “human intelligence tasks (hits)”, such as labeling data. The researcher sets a price for each set of annotations and a list of minimal qualifications for annotators, such as their native language and their satisfaction rate on previous tasks. The use of relatively untrained “crowdworkers” contrasts with earlier annotation efforts, which relied on professional linguists (Marcus et al., 1993). However, crowdsourcing has been found to produce reliable annotations for many language-related tasks (Snow et al., 2008). Crowdsourcing is part of the broader field of **human computation** (Law and Ahn, 2011).

¹¹ For other types of annotations, Krippendorff’s alpha is a popular choice (Hayes and Krippendorff, 2007; Artstein and Poesio, 2008).

Exercises

1. As noted in § 3.3.3, words tend to appear in clumps, with subsequent occurrences of a word being more probable. More concretely, if word j has probability $\phi_{y,j}$ of appearing in a document with label y , then the probability of two appearances ($x_j^{(i)} = 2$) is greater than $\phi_{y,j}^2$.

Suppose you are applying Naïve Bayes to a binary classification. Focus on a word j which is more probable under label $y = 1$, so that,

$$\Pr(w = j \mid y = 1) > \Pr(w = j \mid y = 0). \quad [3.13]$$

Now suppose that $x_j^{(i)} > 1$. All else equal, will the classifier overestimate or underestimate the posterior $\Pr(y = 1 \mid \mathbf{x})$?

2. Prove that F-measure is never greater than the arithmetic mean of recall and precision, $\frac{r+p}{2}$. Your solution should also show that F-measure is equal to $\frac{r+p}{2}$ iff $r = p$.
3. Given a binary classification problem in which the probability of the “positive” label is equal to α , what is the expected F -measure of a random classifier which ignores the data, and selects $\hat{y} = +1$ with probability $\frac{1}{2}$? (Assume that $p(\hat{y}) \perp p(y)$.) What is the expected F -measure of a classifier that selects $\hat{y} = +1$ with probability α (also independent of $y^{(i)}$)? Depending on α , which random classifier will score better?
4. Suppose that binary classifiers c_1 and c_2 disagree on $N = 30$ cases, and that c_1 is correct in $k = 10$ of those cases.
 - Write a program that uses primitive functions such as `exp` and `factorial` to compute the **two-tailed** p -value — you may use an implementation of the “choose” function if one is available. Verify your code against the output of a library for computing the binomial test or the binomial CDF, such as `scipy.stats.binom` in Python.
 - Then use a randomized test to try to obtain the same p -value. In each sample, draw from a binomial distribution with $N = 30$ and $\theta = \frac{1}{2}$. Count the fraction of samples in which $k \leq 10$. This is the one-tailed p -value; double this to compute the two-tailed p -value.
 - Try this with varying numbers of bootstrap samples: $M \in \{100, 1000, 5000, 10000\}$. For $M = 100$ and $M = 1000$, run the test 10 times, and plot the resulting p -values.
 - Finally, perform the same tests for $N = 70$ and $k = 25$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

5. SemCor 3.0 is a labeled dataset for word sense disambiguation. You can download it,¹² or access it in `nltk.corpora.semcor`.

Choose a word that appears at least ten times in SemCor (*find*), and annotate its WordNet senses across ten randomly-selected examples, without looking at the ground truth. Use online WordNet to understand the definition of each of the senses.¹³ Have a partner do the same annotations, and compute the raw rate of agreement, expected chance rate of agreement, and Cohen's kappa.

6. Download the Pang and Lee movie review data, currently available from <http://www.cs.cornell.edu/people/pabo/movie-review-data/>. Hold out a randomly-selected 400 reviews as a test set.

Download a sentiment lexicon, such as the one currently available from Bing Liu, <https://www.cs.uic.edu/~liub/FBS/sentiment-analysis.html>. Tokenize the data, and classify each document as positive iff it has more positive sentiment words than negative sentiment words. Compute the accuracy and F -measure on detecting positive reviews on the test set, using this lexicon-based classifier.

Then train a discriminative classifier (averaged perceptron or logistic regression) on the training set, and compute its accuracy and F -measure on the test set.

Determine whether the differences are statistically significant, using two-tailed hypothesis tests: Binomial for the difference in accuracy, and bootstrap for the difference in macro- F -measure.

The remaining problems will require you to build a classifier and test its properties. Pick a multi-class text classification dataset, such as RCV1¹⁴). Divide your data into training (60%), development (20%), and test sets (20%), if no such division already exists.

7. Compare various vocabulary sizes of 10^2 , 10^3 , 10^4 , 10^5 , using the most frequent words in each case (you may use any reasonable tokenizer). Train logistic regression classifiers for each vocabulary size, and apply them to the development set. Plot the accuracy and Macro- F -measure with the increasing vocabulary size. For each vocabulary size, tune the regularizer to maximize accuracy on a subset of data that is held out from the training set.
8. Compare the following tokenization algorithms:
 - Whitespace, using a regular expression

¹²e.g., https://github.com/google-research-datasets/word_sense_disambiguation_corpus or <http://globalwordnet.org/wordnet-annotated-corpora/>

¹³<http://wordnetweb.princeton.edu/perl/webwn>

¹⁴http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_README.htm

- Penn Treebank
- Split input into five-character units, regardless of whitespace or punctuation

Compute the token/type ratio for each tokenizer on the training data, and explain what you find. Train your classifier on each tokenized dataset, tuning the regularizer on a subset of data that is held out from the training data. Tokenize the development set, and report accuracy and Macro- F -measure.

9. Apply the Porter and Lancaster stemmers to the training set, using any reasonable tokenizer, and compute the token/type ratios. Train your classifier on the stemmed data, and compute the accuracy and Macro- F -measure on stemmed development data, again using a held-out portion of the training data to tune the regularizer.
10. Identify the best combination of vocabulary filtering, tokenization, and stemming from the previous three problems. Apply this preprocessing to the test set, and compute the test set accuracy and Macro- F -measure. Compare against a baseline system that applies no vocabulary filtering, whitespace tokenization, and no stemming.

Use the binomial test to determine whether your best-performing system is significantly more accurate than the baseline.

Use the bootstrap test with $M = 10^4$ to determine whether your best-performing system achieves significantly higher macro- F -measure.

Chapter 4

Learning without supervision

So far we've assumed the following setup:

- A **training set** where you get observations $\mathbf{x}^{(i)}$ and labels $y^{(i)}$
- A **test set** where you only get observations $\mathbf{x}^{(i)}$

If you never get any labeled instances, is it possible to learn anything? This scenario is known as **unsupervised learning**, and we will see that indeed it is possible to learn about the underlying structure of unlabeled observations. We will also explore some related scenarios: **semi-supervised learning**, in which only some instances are labeled, and **domain adaptation**, in which the training data differs from the data on which the trained system will be deployed.

4.1 Unsupervised learning

To motivate unsupervised learning, consider the problem of word sense disambiguation (§ 3.2). Our goal is to classify each instance of a word, such as *bank* into a sense,

- bank#1: a financial institution
- bank#2: the land bordering a river

It is difficult to obtain sufficient training data for word sense disambiguation, because even a large corpus will contain only a few instances of all but the most common words. Unsupervised learning — **word sense induction** — would be highly desirable for this problem.

Word sense disambiguation is usually performed using feature vectors constructed from the local context of the word to be disambiguated. For example, for the word

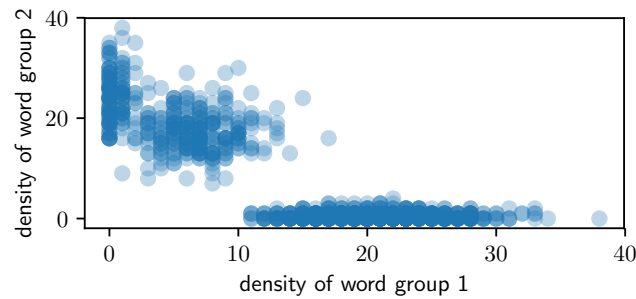


Figure 4.1: Counts of words from two different context groups

bank, the immediate context might typically include words from one of the following two groups:

1. *deposits, interest, regulation*
2. *sediment, river, sand*

Suppose we were to scatterplot each instance of *bank* 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 plot, shown in Figure 4.1, two or more “blobs” might emerge, and these blobs correspond to the different senses of *bank*. These blobs represent the underlying structure of the unlabeled data.

While this points towards a solution, we don’t know the groupings of context words in advance. We have to apply the same idea without predefined word groupings, by working in a high-dimensional space, with one dimension for every context word. Although we can’t plot this space on paper, the idea is the same: try to identify the underlying structure of x , such that there are a few clusters of points, each of which is internally coherent.

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 group of words (*hurricane, winds, storm*); the y -axis is the frequency of the group (*election, voters, vote*). This time, 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 show that even without labels, we can often find structure in data. In two dimensions, groupings of data points can be identified by eye, but this is not possible in the high dimensional scenarios that we face in natural language processing. Instead, we will employ **clustering** algorithms to find this structure automatically.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 8 K -means clustering algorithm

```

1: procedure  $K$ -MEANS( $\mathbf{x}_{1:N}, K$ )
2:   for  $i \in 1 \dots N$  do                                     ▷ initialize cluster memberships
3:      $z^{(i)} \leftarrow \text{RandomInt}(1, K)$ 
4:   repeat
5:     for  $k \in 1 \dots K$  do                                     ▷ recompute cluster centers
6:        $\boldsymbol{\nu}_k \leftarrow \frac{1}{\delta(z^{(i)}=k)} \sum_{i=1}^N \delta(z^{(i)} = k) \mathbf{x}^{(i)}$ 
7:     for  $i \in 1 \dots N$  do                                     ▷ reassign instances to nearest clusters
8:        $z^{(i)} \leftarrow \text{argmin}_k \|\mathbf{x}^{(i)} - \boldsymbol{\nu}_k\|^2$ 
9:   until converged
10:  return  $\{z^{(i)}\}$                                            ▷ return cluster assignments

```

4.1.1 K -means clustering

Clustering algorithms assign each data point to a discrete cluster, $z_i \in 1, 2, \dots, K$. A classical clustering algorithm is **K -means**, an iterative algorithm that maintains a cluster assignment for each instance and a central location for each cluster. K -means iterates between updates to the assignments and the centers:

1. each instance is placed in the cluster with the closest center;
2. each center is recomputed as the average over points in the cluster.

This is formalized in Algorithm 8. The term $\|\mathbf{x}^{(i)} - \boldsymbol{\nu}\|^2$ refers to the squared Euclidean norm, $\sum_n (x_n^{(i)} - \nu_n)^2$.

Soft K -means is a particularly relevant variant. Instead of directly assigning each point to a specific cluster, soft K -means assigns each point a **distribution** over clusters $q^{(i)}$, so that $\sum_{k=1}^K q^{(i)}(k) = 1$, and $\forall_k, 0 \leq q^{(i)}(k) \leq 1$. The soft weight $q^{(i)}(k)$ is computed from the distance of $\mathbf{x}^{(i)}$ to the cluster center $\boldsymbol{\nu}_k$. In turn, the center of each cluster is computed from a **weighted average** of the points in the cluster,

$$\boldsymbol{\nu}_k = \frac{1}{\sum_{i=1}^N q^{(i)}(k)} \sum_{i=1}^N q^{(i)}(k) \mathbf{x}^{(i)}. \quad [4.1]$$

We will now explore a probabilistic version of soft K -means clustering, based on **expectation maximization** (EM). Because EM clustering can be derived as an approximation to maximum-likelihood estimation, it can be extended in a number of useful ways.

(c) Jacob Eisenstein 2014-2017. Work in progress.

4.1.2 Expectation Maximization (EM)

Expectation maximization combines the idea of soft K -means with Naïve Bayes classification. To review, Naïve Bayes defines a probability distribution over the data,

$$\log p(\mathbf{x}, \mathbf{y}; \phi, \mu) = \sum_{i=1}^N \log \left(p(\mathbf{x}^{(i)} | y^{(i)}; \phi) \times p(y^{(i)}; \mu) \right) \quad [4.2]$$

Suppose we never observe the labels. To indicate this, we'll refer to the label of each instance as $z^{(i)}$, rather than $y^{(i)}$, which is usually reserved for observed variables. By marginalizing over the **latent** variables z , we compute a distribution over the observed instances \mathbf{x} :

$$\log p(\mathbf{x}; \phi, \mu) = \sum_{i=1}^N \log p(\mathbf{x}^{(i)}; \phi, \mu) \quad [4.3]$$

$$= \sum_{i=1}^N \log \sum_{z=1}^K p(\mathbf{x}^{(i)}, z; \phi, \mu) \quad [4.4]$$

$$= \sum_{i=1}^N \log \sum_{z=1}^K p(\mathbf{x}^{(i)} | z; \phi) \times p(z; \mu). \quad [4.5]$$

To estimate the parameters ϕ and μ , we can maximize the marginal likelihood in Equation 4.5. Why is this the right 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.

When the labels are observed, we can estimate the parameters of the Naïve Bayes probability model separately for each label. But marginalizing over the labels couples these parameters, making direct optimization of $\log p(\mathbf{x})$ intractable. We will approximate the log-likelihood by introducing an **auxiliary variable** $q^{(i)}$, which is a distribution over the label set $\mathcal{Z} = \{1, 2, \dots, K\}$. Our optimization procedure will alternate between updates to q and updates to the parameters (ϕ, μ) . Thus, $q^{(i)}$ plays the same role in EM as in soft K -means.

To derive the updates for this optimization, multiply the right side of Equation 4.5 by

(c) Jacob Eisenstein 2014-2017. Work in progress.

the ratio $\frac{q^{(i)}(z)}{q^{(i)}(z)} = 1$,

$$\log p(\mathbf{x}; \phi, \mu) = \sum_{i=1}^M \log \sum_{z=1}^K p(\mathbf{x}^{(i)} | z; \phi) \times p(z; \mu) \times \frac{q^{(i)}(z)}{q^{(i)}(z)} \quad [4.6]$$

$$= \sum_{i=1}^M \log \sum_{z=1}^K q^{(i)}(z) \times p(\mathbf{x}^{(i)} | z; \phi) \times p(z; \mu) \times \frac{1}{q^{(i)}(z)} \quad [4.7]$$

$$= \sum_{i=1}^M \log E_{q^{(i)}} \left[\frac{p(\mathbf{x}^{(i)} | z; \phi) p(z; \mu)}{q^{(i)}(z)} \right], \quad [4.8]$$

where $E_{q^{(i)}} [f(z)] = \sum_{z=1}^K q^{(i)}(z) \times f(z)$ refers to the expectation of the function f under the distribution $z \sim q^{(i)}$.

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

$$\log p(\mathbf{x}; \phi, \mu) \geq \sum_{i=1}^N E_{q^{(i)}} \left[\log \frac{p(\mathbf{x}^{(i)} | z; \phi) p(z; \mu)}{q^{(i)}(z)} \right] \quad [4.9]$$

$$\mathcal{J} \triangleq \sum_{i=1}^N E_{q^{(i)}} \left[\log p(\mathbf{x}^{(i)} | z; \phi) + \log p(z; \mu) - \log q^{(i)}(z) \right] \quad [4.10]$$

$$= \sum_{i=1}^N E_{q^{(i)}} \left[\log p(\mathbf{x}^{(i)}, z; \phi, \mu) \right] + H(q^{(i)}) \quad [4.11]$$

We will focus on Equation 4.10, which is the lower bound on the marginal log-likelihood of the observed data, $\log p(\mathbf{x})$. Equation 4.11 shows the connection to the information theoretic concept of **entropy**, $H(q^{(i)}) = -\sum_z q^{(i)}(z) \log q^{(i)}(z)$, which measures the average amount of information produced by a draw from the distribution $q^{(i)}$. Returning to Equation 4.10, this bound is a function of two groups of arguments:

- the distributions $q^{(i)}$ for each instance;
- the parameters μ and ϕ .

We will optimize with respect to each of these in turn, while holding the other fixed.

The E-step

The step in which we update $q^{(i)}$ is known as the **E-step**, because we are updating the distribution under which the expectation is computed. To derive this update, we first

write out the expectation in the lower bound as a sum,

$$\mathcal{J} = \sum_{i=1}^N \sum_{z=1}^K q^{(i)}(z) \left[\log p(\mathbf{x}^{(i)} | z; \phi) + \log p(z; \mu) - \log q^{(i)}(z) \right]. \quad [4.12]$$

As in Naïve Bayes, we have a “sum-to-one” constraint: in this case, $\sum_{z=1}^K q^{(i)}(z) = 1$. Once again, we incorporate this constraint into a Lagrangian:

$$\mathcal{J}_q = \sum_{i=1}^N \sum_{z=1}^K q^{(i)}(z) \left(\log p(\mathbf{x}^{(i)} | z; \phi) + \log p(z; \mu) - \log q^{(i)}(z) \right) + \lambda^{(i)} \left(1 - \sum_{z=1}^K q^{(i)}(z) \right) \quad [4.13]$$

We then optimize by taking the derivative and setting it equal to zero:

$$\frac{\partial \mathcal{J}_q}{\partial q^{(i)}(z)} = \log p(\mathbf{x}^{(i)} | z; \phi) + \log p(z; \mu) - \log q^{(i)}(z) - 1 - \lambda^{(i)} \quad [4.14]$$

$$\log q^{(i)}(z) = \log p(\mathbf{x}^{(i)} | z; \phi) + \log p(z; \mu) - 1 - \lambda^{(i)} \quad [4.15]$$

$$q^{(i)}(z) \propto p(\mathbf{x}^{(i)} | z; \phi) \times p(z; \mu). \quad [4.16]$$

Since $q^{(i)}$ is constrained to be a probability distribution, the normalized probability can be computed easily,

$$q^{(i)}(z) = \frac{p(\mathbf{x}^{(i)} | z; \phi) \times p(z; \mu)}{\sum_{z'=1}^K p(\mathbf{x}^{(i)} | z'; \phi) \times p(z'; \mu)} \quad [4.17]$$

$$= p(z | \mathbf{x}^{(i)}; \phi, \mu). \quad [4.18]$$

After normalizing, each $q^{(i)}$ — which is the soft distribution over clusters for data $\mathbf{x}^{(i)}$ — is set to the posterior probability $p(z | \mathbf{x}^{(i)}; \phi, \mu)$ under the current parameters. Although the Lagrange multipliers $\lambda^{(i)}$ were introduced as additional parameters, they drop out during normalization.

The M-step

Next, we hold fixed the soft assignments $q^{(i)}$, and maximize with respect to the parameters, ϕ and μ . We’ll focus on the parameter ϕ , which parametrizes the likelihood, $p(\mathbf{x} | z; \phi)$, and leave μ for an exercise. The parameter ϕ is a distribution over words for each cluster, so it is optimized under constraint that $\sum_{j=1}^V \phi_{z,j} = 1$. To incorporate this constraint, we again form a Lagrangian,

$$\mathcal{J}_\phi = \sum_{i=1}^N \sum_{z=1}^K q^{(i)}(z) \left(\log p(\mathbf{x}^{(i)} | z; \phi) + \log p(z; \mu) - \log q^{(i)}(z) \right) + \sum_{z=1}^K \lambda_z \left(1 - \sum_{j=1}^V \phi_{z,j} \right). \quad [4.19]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

The conditional log-likelihood for the multinomial is,

$$\log p(\mathbf{x}^{(i)} \mid z, \phi) = C + \sum_{j=1}^V x_j \log \phi_{z,j}, \quad [4.20]$$

where C is a constant with respect to ϕ . Setting the derivative of \mathcal{J}_ϕ equal to zero,

$$\frac{\partial \mathcal{J}_\phi}{\partial \phi_{z,j}} = \sum_{i=1}^N q^{(i)}(z) \times \frac{x_j^{(i)}}{\phi_{z,j}} - \lambda_z \quad [4.21]$$

$$\phi_{z,j} \propto \sum_{i=1}^N q^{(i)}(z) \times x_j^{(i)}. \quad [4.22]$$

Because ϕ_z is constrained to be a probability distribution, the normalized probability can be computed easily,

$$\phi_{z,j} = \frac{\sum_{i=1}^N q^{(i)}(z) \times x_j^{(i)}}{\sum_{j'=1}^V \sum_{i=1}^N q^{(i)}(z) \times x_{j'}^{(i)}} = \frac{E_q [\text{count}(z, j)]}{\sum_{j'=1}^V E_q [\text{count}(z, j')]}, \quad [4.23]$$

where the counter $j \in \{1, 2, \dots, V\}$ indexes over base features, such as words.

This update sets ϕ_z equal to the relative frequency estimate of the **expected counts** under the distribution q . As in supervised Naïve Bayes, we can smooth these counts by adding a constant α . The update for μ is similar: $\mu_z \propto \sum_{i=1}^N q^{(i)}(z) = E_q [\text{count}(z)]$, the expected frequency of cluster z . This probability can also be smoothed. In sum, the M-step is just like Naïve Bayes, but with expected counts rather than observed counts.

The multinomial likelihood $p(\mathbf{x} \mid z)$ can be replaced with other probability distributions: for example, for continuous observations, a Gaussian distribution can be used. In some cases, there is no closed-form update to the parameters of the likelihood. One approach is to run gradient-based optimization at each M-step; another is to simply take a single step along the gradient step and then return to the E-step (Berg-Kirkpatrick et al., 2010).

4.1.3 EM as an optimization algorithm

Algorithms that alternate between updating various subsets of the parameters are called **coordinate ascent** algorithms. The objective function \mathcal{J} is **biconvex**: it is separately convex in q and $\langle \mu, \phi \rangle$, but it is not jointly convex in all terms. In the coordinate ascent algorithm that we have defined, each step is guaranteed not to decrease \mathcal{J} . Thanks to this guarantee, we know that EM will converge towards a solution at which no nearby points yield further improvements. This solution is a **local optimum** — it is as good or

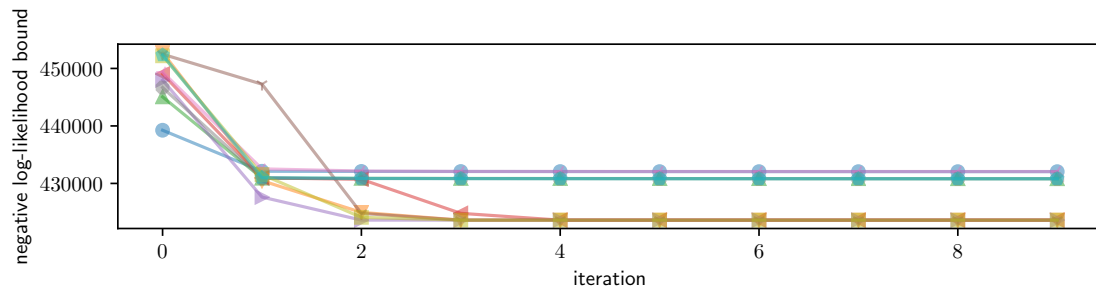


Figure 4.2: Sensitivity of expectation maximization to initialization. Each line shows the progress of optimization from a different random initialization.

better than any of its immediate neighbors, but is **not** guaranteed to be optimal among all possible configurations of $\langle \mathbf{q}, \boldsymbol{\mu}, \boldsymbol{\phi} \rangle$.

The fact that there is no guarantee of global optimality means that initialization is important: where you start can determine where you finish. To illustrate this point, Figure 4.2 shows the objective function for EM with ten different random initializations: while the objective function improves monotonically in each run, it converges to several different values.¹ For linear classifiers like logistic regression, we don't need to worry about initialization, because it won't affect the ultimate solution: we are guaranteed to reach the global minimum. (Initialization does matter for more complex supervised learning algorithms like multilayer neural networks.) Recent work on **spectral learning** has sought to obtain similar guarantees for unsupervised learning; this is briefly described in § 4.5.

In **hard EM**, each $\mathbf{q}^{(i)}$ distribution assigns probability of 1 to a single label $\hat{z}^{(i)}$, and zero probability to all others (Neal and Hinton, 1998). This is similar in spirit to K -means clustering, and can outperform standard EM in some cases (Spitkovsky et al., 2010). Another variant of coordinate ascent combines EM with stochastic gradient descent (SGD). In this case, we can do a local E-step at each instance $\mathbf{x}^{(i)}$, and then immediately make a gradient update to the parameters $\langle \boldsymbol{\mu}, \boldsymbol{\phi} \rangle$. This algorithm has been called **incremental expectation maximization** (Neal and Hinton, 1998) and **online expectation maximization** (Sato and Ishii, 2000; Cappé and Moulines, 2009), and is especially useful when there is no closed-form optimum for the likelihood $p(\mathbf{x} | \mathbf{z})$, and in online settings where new data is constantly streamed in (see Liang and Klein, 2009, for a comparison for online EM variants).

¹The figure shows the bound on the **negative** log-likelihood, because optimization is typically framed as minimization rather than maximization.

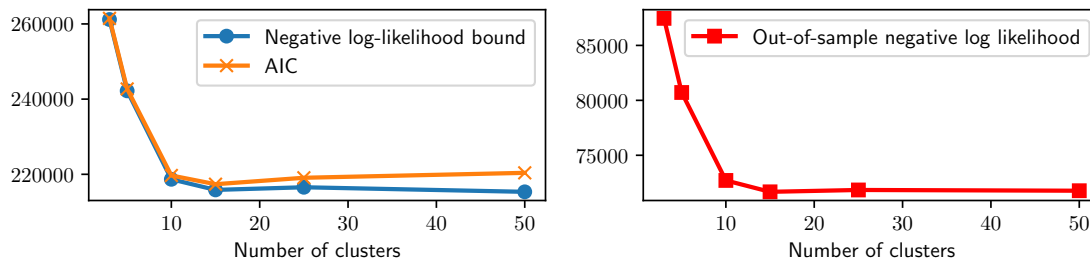


Figure 4.3: The negative log-likelihood and AIC for several runs of expectation maximization, on synthetic data. Although the data was generated from a model with $K = 10$, the optimal number of clusters is $\hat{K} = 15$, according to AIC and the heldout log-likelihood. The training set log-likelihood continues to improve as K increases.

4.1.4 How many clusters?

So far, we have assumed that the number of clusters K is given. In some cases, this assumption is valid. For example, a lexical semantic resource like WordNet might tell us the number of senses for a word. The number of clusters might also be a parameter for the user to tune: some readers want a coarse-grained clustering of news stories into three or four clusters, while others want a fine-grained clustering into twenty or more. But in many cases, we must choose K automatically, with little extrinsic guidance.

One solution is to choose the number of clusters to maximize metric of clustering quality. The other parameters μ and ϕ are chosen to maximize the log-likelihood bound \mathcal{J} , so this might seem a potential candidate for tuning K . However, \mathcal{J} will never decrease with K : if it is possible to obtain a bound of \mathcal{J}_K with K clusters, then it is always possible to do at least as well with $K + 1$ clusters, by simply ignoring the additional cluster and setting its probability to zero in \mathbf{q} and μ . It is therefore necessary to introduce a penalty for model complexity, so that fewer clusters are preferred. For example, the Akaike Information Criterion (AIC; Akaike, 1974) is the linear combination of the number of parameters and the log-likelihood,

$$\text{AIC} = 2M - 2\mathcal{J}, \quad [4.24]$$

where M is the number of parameters. In an expectation maximization clustering algorithm, $M = K \times V + K$. 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 , and so it is not the case that the likelihood on this data is guaranteed to increase with K . Figure 4.3 shows the negative log-likelihood on training and heldout data, as well as the AIC.

***Bayesian nonparametrics** An alternative approach is to treat the number of clusters as another latent variable. This requires statistical inference over a set of models with a variable number of clusters. This is not possible within the framework of expectation maximization, but there are several alternative inference procedures which can be applied, including **Markov Chain Monte Carlo (MCMC)**, which is briefly discussed in § 4.5 (for more details, see Chapter 25 of Murphy, 2012). Bayesian nonparametrics have been applied to the problem of unsupervised word sense induction, thus learning not only the word senses but also the number of senses per word (Reisinger and Mooney, 2010).

4.2 Applications of expectation-maximization

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 z , such that it is easy to write the probability $P(\mathcal{D}, z)$, where \mathcal{D} is the observed data. It should also be easy to estimate the associated parameters, given knowledge of z .
- Derive the E-step updates for $q(z)$, which is typically factored as $q(z) = \prod_{i=1}^N q_{z^{(i)}}(z^{(i)})$, where i is an index over instances.
- The M-step updates typically correspond to the soft version of a probabilistic supervised learning algorithm, like Naïve Bayes.

Here are a few of the many applications of this general framework.

4.2.1 Word sense induction

The chapter began by considering the problem of word sense disambiguation when the senses are not known in advance. Expectation-maximization can be applied to this problem by treating each cluster as a word sense. Each instance represents the use of an ambiguous word, and $\mathbf{x}^{(i)}$ is a vector of counts for the other words that appear nearby: Schütze (1998) uses all words within a 50-word window. The probability $p(\mathbf{x}^{(i)} | y)$ can be set to the multinomial distribution, as in Naïve Bayes. The EM algorithm can be applied directly to this data, yielding clusters that (hopefully) correspond to the word senses.

Better performance can be obtained by first applying truncated **singular value decomposition (SVD)** to the matrix of context-counts $\mathbf{C}_{ij} = \text{count}(i, j)$, where $\text{count}(i, j)$ is the

(c) Jacob Eisenstein 2014-2017. Work in progress.

count of word j in the context of instance i .

$$\begin{aligned}
& \min_{\mathbf{U}, \mathbf{S}, \mathbf{V}} \|\mathbf{C} - \mathbf{USV}^\top\|_F & [4.25] \\
& s.t. \mathbf{U} \in \mathbb{R}^{V \times K}, \mathbf{UU}^\top = \mathbb{I} \\
& \mathbf{S} = \text{Diag}(s_1, s_2, \dots, s_K) \\
& \mathbf{V}^\top \in \mathbb{R}^{N_p \times K}, \mathbf{VV}^\top = \mathbb{I},
\end{aligned}$$

where $\|\cdot\|_F$ is the Frobenius norm, $\|X\|_F = \sqrt{\sum_{i,j} X_{i,j}^2}$. The matrix \mathbf{U} contains the left singular vectors of \mathbf{C} , and the rows of this matrix can be used as low-dimensional representations of the count vectors \mathbf{c}_i . EM clustering can be made more robust by setting the instance descriptions $\mathbf{x}^{(i)}$ equal to these rows, rather than using raw counts (Schütze, 1998). However, because the instances are now dense vectors of continuous numbers, the probability $p(\mathbf{x}^{(i)} | z)$ must be defined as a multivariate Gaussian distribution.

In truncated singular value decomposition, the hyperparameter K is the truncation limit: when K is equal to the rank of \mathbf{C} , the norm of the difference between the original matrix \mathbf{C} and its reconstruction \mathbf{USV}^\top will be zero. Lower values of K increase the reconstruction error, but yield vector representations that are smaller and easier to learn from. Singular value decomposition is discussed in more detail in chapter 13.

4.2.2 Semi-supervised learning

Expectation-maximization can also be applied to the problem of **semi-supervised learning**: learning from both labeled and unlabeled data in a single model. Semi-supervised learning makes use of ground truth annotations, ensuring that each label y corresponds to the desired concept. At the same time, by making use of unlabeled data, it is possible cover a greater fraction of the features than would be possible using labeled data alone. Other methods for semi-supervised learning are discussed in § 4.3, but for now, we approach the problem within the framework of expectation-maximization (Nigam et al., 2000).

Suppose we have labeled data $\{\mathbf{x}^{(i)}, y^{(i)}\}_{i=1}^{N_\ell}$, and unlabeled data $\{\mathbf{x}^{(i)}\}_{i=N_\ell+1}^{N_\ell+N_u}$, where N_ℓ is the number of labeled instances and N_u is the number of unlabeled instances. We can learn from the combined data by maximizing the joint log-likelihood,

$$\mathcal{L} = \sum_{i=1}^{N_\ell} \log p(\mathbf{x}^{(i)}, y^{(i)}; \boldsymbol{\mu}, \phi) + \sum_{j=N_\ell+1}^{N_\ell+N_u} \log p(\mathbf{x}^{(j)}; \boldsymbol{\mu}, \phi) \quad [4.26]$$

$$= \sum_{i=1}^{N_\ell} \left(\log p(\mathbf{x}^{(i)} | y^{(i)}; \phi) + \log p(y^{(i)}; \boldsymbol{\mu}) \right) + \sum_{j=N_\ell+1}^{N_\ell+N_u} \log \sum_{y=1}^K p(\mathbf{x}^{(j)}, y; \boldsymbol{\mu}, \phi). \quad [4.27]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

The left sum is identical to the objective in Naïve Bayes; the right sum is the marginal log-likelihood for expectation-maximization clustering, from Equation 4.5. We can construct a lower bound on this log-likelihood by introducing distributions $q^{(j)}$ for all $j \in \{N_\ell + 1, \dots, N_\ell + N_u\}$. During the E-step, we update these distributions; during the M-step we update the parameters ϕ and μ using the expected counts from the unlabeled data and the observed counts from the labeled data.

A critical issue in semi-supervised learning is how to balance the impact of the labeled and unlabeled data on the classifier weights, especially when the unlabeled data is much larger than the labeled dataset. The risk is that the unlabeled data will dominate, causing the parameters to drift towards a “natural clustering” of the instances — which may not correspond to a good classifier for the labeled data. One solution is to heuristically reweight the two components of Equation 4.26, which can be critical to achieving good performance with EM (Nigam et al., 2000).

4.2.3 Multi-component modeling

As a final application, let’s return to fully 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; if so, it would be better to model these components separately, rather than treating their union as a single class. However, the component responsible for each instance is not directly observed.

Recall that Naïve Bayes is based on a generative process, which provides a stochastic explanation for the observed data. In Naïve Bayes, each label is drawn from a categorical distribution with parameter μ , and each vector of word counts is drawn from a multinomial distribution with parameter ϕ_y . 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

(c) Jacob Eisenstein 2014-2017. Work in progress.

-
- (4.1) ☹ Villeneuve a bel et bien **réussi** son pari de changer de perspectives tout en assurant une cohérence à la franchise.²
- (4.2) ☹ Il est également trop **long** et bancal dans sa narration, tiède dans ses intentions, et tiraillé entre deux personnages et directions qui ne parviennent pas à coexister en harmonie.³
- (4.3) Denis Villeneuve a **réussi** une suite **parfaitement** maîtrisée⁴
- (4.4) **Long**, **bavard**, hyper design, à peine agité (le comble de l'action : une bagarre dans la flotte), métaphysique et, surtout, ennuyeux jusqu'à la catalepsie.⁵
- (4.5) Une suite d'une écrasante puissance, mêlant **parfaitement** le contemplatif au narratif.⁶
- (4.6) Le film impitoyablement **bavard** finit quand même par se taire quand se lève l'espèce de bouquet final où semble se déchaîner, comme en libre parcours de poulets décapités, l'armée des graphistes numériques griffant nerveusement la palette graphique entre agonie et orgasme.⁷
-

Table 4.1: Labeled and unlabeled reviews of the films *Blade Runner 2049* and *Transformers: The Last Knight*.

lower bound on the log-likelihood,

$$\log p(\mathbf{x}^{(i)}, y^{(i)}) = \log \sum_{z=1}^{K_z} p(\mathbf{x}^{(i)}, y^{(i)}, z; \boldsymbol{\mu}, \boldsymbol{\phi}, \boldsymbol{\beta}) \quad [4.28]$$

$$\geq \log p(y^{(i)}; \boldsymbol{\mu}) + E_q \left[\log p(\mathbf{x}^{(i)} | z; \boldsymbol{\phi}) + \log p(z | y^{(i)}; \boldsymbol{\beta}) - \log q^{(i)}(z) \right]. \quad [4.29]$$

We are now ready to apply expectation maximization. As usual, E-step updates the distribution over the missing data, $q^{(i)}$. During the M-step, we update the parameters,

$$\beta_{y,z} = \frac{E_q [\text{count}(y, z)]}{\sum_{z'=1}^{K_z} E_q [\text{count}(y, z')]} \quad [4.30]$$

$$\phi_{z,j} = \frac{E_q [\text{count}(z, j)]}{\sum_{j'=1}^V E_q [\text{count}(z, j')]} \quad [4.31]$$

4.3 Semi-supervised learning

In semi-supervised learning, the learner makes use of both labeled and unlabeled data. To see how this could help, suppose you want to do sentiment analysis in French. In Table 4.1, there are two labeled examples, one positive and one negative. From this data, a learner conclude that *réussi* is positive and *long* is negative; this isn't much! However, we can propagate this information to the unlabeled data, and potentially learn more.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- If we are confident that *réussi* is positive, then we might guess that (4.3) is also positive.
- That suggests that *parfaitement* is also positive.
- We can then propagate this information to (4.5), and learn more.
- Similarly, we can propagate from the labeled data to (4.4), which we guess to be negative because it shares the word *long*. This suggests that *bavard* is also negative, which we propagate to (4.6).

Instances (4.3) and (4.4) were “similar” to our labeled examples for positivity and negativity, respectively. By using these instances to expand the models for each class, it became possible to correctly label instances (4.5) and (4.6), which didn’t share any important features with our original labeled data. This requires a key assumption: that similar instances will have similar labels.

In § 4.2.2, we saw how expectation-maximization can be applied to semi-supervised learning. In this case, you can imagine that the initial parameters ϕ assigned a high weight for *réussi* in the positive class, and a high weight for *long* in the negative class. These weights helped to shape the distributions q for instances (4.3) and (4.4) in the E-step. In the next iteration of the M-step, the parameters ϕ are updated with counts from these instances, making it possible to correctly label the instances (4.5) and (4.6).

However, expectation-maximization has a key disadvantage: it requires using a generative classification model, which imposes severe restrictions on the features that can be used for classification. In this section, we explore non-probabilistic approaches, which impose fewer restrictions on the classification model.

4.3.1 Multi-view learning

EM semi-supervised learning can be viewed as **self-training**: the labeled data guides the initial estimates of the classification parameters; these parameters are used to compute a label distribution over the unlabeled instances, $q^{(i)}$; the label distributions are used to update the parameters. The risk is that self-training drifts away from the original labeled data. This problem can be ameliorated by multi-view learning. The key assumption is that the features can be decomposed into multiple “views”, each of which is conditionally independent, given the label. For example, consider the problem of classifying a name as a person or location: one view is the name itself; another is the context in which it appears. This situation is illustrated in Table 4.2.

Co-training is an iterative multi-view learning algorithm, in which we have separate classifiers for each view (Blum and Mitchell, 1998). At each iteration of the algorithm, each classifier predicts labels for a subset of the unlabeled instances, using only the features available in its view. These predictions are then used as ground truth to train the classifiers

	$x^{(1)}$	$x^{(2)}$	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

Table 4.2: Example of multiview learning for named entity classification

associated with the other views. In the example shown in Table 4.2, the classifier on $x^{(1)}$ might correctly label instance #5 as a person, because of the feature *Dr*; this instance would then serve as training data for the classifier on $x^{(2)}$, which would then be able to correctly label instance #6, thanks to the feature *recommended*. If the views are truly independent, this procedure is robust to drift. Furthermore, it imposes no restrictions on the classifiers that can be used for each view.

Word-sense disambiguation is particularly suited to multi-view learning, thanks to the heuristic of “one sense per discourse”: if a polysemous word is used more than once in a given text or conversation, all usages refer to the same sense (Gale et al., 1992). This motivates a multi-view learning approach, in which one view corresponds to the local context, and another view corresponds to the global context at the document level (Yarowsky, 1995). We first train the local context view on a small seed dataset, and then identify its most confident predictions on unlabeled instances. The global context view is then used to extend these confident predictions to other instances within the same document. These new instances are added to the training data to the local context classifier, which is retrained and then re-applied to the remaining unlabeled data.

4.3.2 Graph-based algorithms

Another family of approaches to semi-supervised learning begins by constructing a graph, in which pairs of instances are linked with symmetric weights $w_{i,j}$, e.g.,

$$w_{i,j} = \exp(-\alpha \times \|\mathbf{x}^{(i)} - \mathbf{x}^{(j)}\|^2). \quad [4.32]$$

Our goal is to use this weighted graph to propagate labels from a small set of labeled instances to larger set of unlabeled instances.

In **label propagation**, this is done through a series of matrix operations (Zhu et al., 2003). Let \mathbf{Q} be a matrix of size $N \times K$, in which each row $\mathbf{q}^{(i)}$ describes the labeling of instance i . When ground truth labels are available, we set $\mathbf{q}^{(i)}$ to an indicator vector, with $q_{y^{(i)}}^{(i)} = 1$ and $q_{y' \neq y^{(i)}}^{(i)} = 0$. We will refer to the submatrix of rows containing labeled

instances as \mathbf{Q}_L , and the remaining rows as \mathbf{Q}_U . The rows of \mathbf{Q}_U can be initialized to assign equal probabilities to all labels, $q_{i,k} = \frac{1}{K}$.

Now, let $T_{i,j}$ represent the “transition” probability of moving from node j to node i ,

$$T_{i,j} \triangleq \Pr(j \rightarrow i) = \frac{w_{i,j}}{\sum_{l=1}^N w_{l,j}}. \quad [4.33]$$

We compute values of $T_{i,j}$ for all instances j and all **unlabeled** instances i , forming a matrix of size $N_U \times N$. If the dataset is large, this matrix may be expensive to store and manipulate; a solution is to sparsify it, by keeping only the κ largest values in each row, and setting all other values to zero. We can then “propagate” the label distributions to the unlabeled instances,

$$\tilde{\mathbf{Q}}_U \leftarrow \mathbf{T}\mathbf{Q} \quad [4.34]$$

$$\mathbf{s} \leftarrow \tilde{\mathbf{Q}}_U \mathbf{1} \quad [4.35]$$

$$\mathbf{Q}_U \leftarrow \text{Diag}(\mathbf{s})^{-1} \tilde{\mathbf{Q}}_U. \quad [4.36]$$

The second two lines normalize the rows of $\tilde{\mathbf{Q}}$, so that each row of \mathbf{Q} is a probability distribution over labels.

4.4 Domain adaptation

We now consider a scenario in which the labeled data differs in some key respect from the unlabeled data to which we want to apply our trained model. A classic example is in consumer reviews: we may have labeled reviews of movies (the **source domain**), but we want to predict the reviews of appliances (the **target domain**). Similar issues arise in cases of genre: most linguistically-annotated data is in the domain of news text, but application domains include social media, electronic health records, and the text of government regulations. In general, there may be several source and target domains, each with their own properties; however, for simplicity we will focus mainly on the case of a single source and target domain.

The simplest approach is “direct transfer”: train a classifier on the source domain, and apply it directly to the target domain. The accuracy of this approach will depend on the extent to which features are shared across domains. In review text, words like *outstanding* and *disappointing* will apply across both movies and appliances; but others, like *terrifying*, may have meanings that are domain-specific. **Domain adaptation** algorithms attempt to do better than direct transfer, by learning from data in both domains. There are two main families of domain adaptation algorithms, depending on whether any labeled data is available in the target domain.

(c) Jacob Eisenstein 2014-2017. Work in progress.

4.4.1 Supervised domain adaptation

In supervised domain adaptation, we assume a small amount of labeled data in the target domain. The simplest approach would be to ignore domain differences, and simply merge the training data from the source and target domains. There are several other baseline approaches to dealing with this scenario (Daumé III, 2007):

- **Interpolation:** train a classifier for each domain, and combine their predictions. For example,

$$\hat{y} = \operatorname{argmax}_y \lambda_s \theta_s \cdot \mathbf{f}_s(\mathbf{x}, y) + (1 - \lambda_s) \theta_t \cdot \mathbf{f}_t(\mathbf{x}, y), \quad [4.37]$$

where θ_s and θ_t are the weights of the source and target classifiers, and λ_s is the interpolation weight on the source domain.

- **Prediction:** train a classifier on the source domain data, use its prediction as an additional feature in a classifier trained on the target domain data.
- **Priors:** train a classifier on the source domain data, and use its weights as a prior distribution on the weights of the classifier for the target domain data. Equivalently, we can regularize the target domain weights towards the weights of the source domain classifier (Chelba and Acero, 2006),

$$\ell(\theta_t) = \sum_{i=1}^N \ell^{(i)}(\mathbf{x}^{(i)}, y^{(i)}; \theta_t) + \lambda \|\theta_t - \theta_s\|_2^2, \quad [4.38]$$

where $\ell^{(i)}$ is the prediction loss on instance i , and λ is the regularization weight.

An effective and “frustratingly simple” alternative is `EasyAdapt` (Daumé III, 2007). The main idea is to create copies of each feature: one for each domain and one for the cross-domain setting. For example, a negative review of the film *Wonder Woman* begins, *As boring and flavorless as a three-day-old grilled cheese sandwich...*⁸ The resulting bag-of-words feature vector would be,

$$\begin{aligned} \mathbf{f}(\mathbf{x}, y, d) = \{ & \langle \text{boring}, -, \text{MOVIE} \rangle, \langle \text{boring}, -, * \rangle, \\ & \langle \text{flavorless}, -, \text{MOVIE} \rangle, \langle \text{flavorless}, -, * \rangle, \\ & \langle \text{three-day-old}, -, \text{MOVIE} \rangle, \langle \text{three-day-old}, -, * \rangle, \\ & \dots \}, \end{aligned}$$

with the asterisk indicating the cross-domain copy of each feature. It is then up to the learner to allocate weight between the domain-specific and cross-domain features: for

⁸<http://www.colesmithey.com/capsules/2017/06/wonder-woman.HTML>, accessed October 9, 2017.

words that facilitate prediction in both domains, the learner will use the cross-domain features; for words that are relevant only to a single domain, the domain-specific features will be used. Any discriminative classifier can be used with these augmented features.⁹

4.4.2 Unsupervised domain adaptation

In unsupervised domain adaptation, we have only unlabeled data in the target domain. Unsupervised domain adaptation algorithms cope with this problem by trying to make the data from the source and target domains as similar as possible. This is typically done by learning a **projection function**, which puts the source and target data in a shared space, in which a learner can generalize across domains. This projection is learned from data in both domains, and is applied to the base features — for example, the bag-of-words in text classification. The projected features can then be used both for training and for prediction.

Linear projection

Many unsupervised domain adaptation algorithms make use of linear projection,

$$\mathbf{f}(\mathbf{x}^{(i)}) = \mathbf{A}\mathbf{x}^{(i)}. \quad [4.39]$$

The projected vectors $\mathbf{f}(\mathbf{x}^{(i)})$ can then be used as base features during both training (from the source domain) and prediction (on the target domain).

The projection matrix \mathbf{U} can be learned in a number of different ways, but many approaches focus on compressing and reconstructing the features in \mathbf{X} (Ando and Zhang, 2005). For example, we can define a set of **pivot features**, which are typically chosen because they appear in both domains: in the case of review documents, pivot features might include evaluative adjectives like *outstanding* and *disappointing* (Blitzer et al., 2007). For each pivot feature j , we define an auxiliary problem of predicting whether the feature is present in each example, using the remaining base features. Let the weights of this classifier be written ϕ_j , and let us horizontally concatenate the weights for each of the N_p pivot features into a matrix $\Phi = [\phi_1, \phi_2, \dots, \phi_{N_p}]$.

We then perform truncated singular value decomposition on Φ , as described in § 4.2.1, obtaining $\Phi \approx \mathbf{U}\mathbf{S}\mathbf{V}^\top$. Now, the rows of the matrix \mathbf{U} summarize information about each base feature: indeed, the truncated singular value decomposition identifies a low-dimension basis for the weight matrix Φ , which in turn links base features to pivot features. Suppose that a base feature *reliable* occurs only in the target domain of appliance reviews. Nonetheless, it will have a positive weight towards some pivot features (e.g., *outstanding*, *recommended*), and a negative weight towards others (e.g., *worthless*, *unpleasant*). A base feature such as *watchable* might have the same associations with the pivot features,

⁹EasyAdapt can be explained as a hierarchical Bayesian model, in which the weights for each domain are drawn from a shared prior (Finkel and Manning, 2009).

and therefore, $\mathbf{u}_{\text{reliable}} \approx \mathbf{u}_{\text{watchable}}$. The matrix \mathbf{U} can thus project the base features into a space in which this information is shared, $\mathbf{f}(\mathbf{x}^{(i)}) = \mathbf{U}^\top \mathbf{x}^{(i)}$.

Non-linear projection

More recent work has focused on learning non-linear transformations of the base features into a shared feature space that enables cross-domain transfer. This is accomplished by implementing the transformation function \mathbf{f} as a deep neural network, which is trained from an auxiliary objective.

Denoising objectives One possibility is to train \mathbf{f} to reconstruct a corrupted version of the original input. The original input can be corrupted in various ways: by the addition of random noise (Glorot et al., 2011; Chen et al., 2012), or by the deletion of features (Chen et al., 2012; Yang and Eisenstein, 2015). Denoising objectives share many properties of the linear projection method described above: they enable \mathbf{f} to be trained on large amounts of unlabeled data from the target domain, and allow information to be shared across the feature space, thereby reducing sensitivity to rare and domain-specific features.

Domain-generality objectives Our ultimate goal is for the transformed representations $\mathbf{f}(\mathbf{x}^{(i)})$ to be domain-general. We can make this an explicit optimization criterion, by computing the similarity of transformed instances both within and between domains (Tzeng et al., 2015), or by formulating an additional classification task, in which the domain itself is treated as a label (Ganin et al., 2016). This setting is **adversarial**, because we want to learn a representation that makes this classifier perform poorly. At the same time, we want $\mathbf{f}(\mathbf{x}^{(i)})$ to enable accurate predictions of the labels $y^{(i)}$.

To formalize this idea, let $d^{(i)}$ represent the domain of instance i , and let $\ell(\mathbf{f}(\mathbf{x}^{(i)}), d^{(i)}; \boldsymbol{\theta}_d)$ represent the loss of a classifier (typically a deep neural network) trained to predict $d^{(i)}$ from the transformed representation $\mathbf{f}(\mathbf{x}^{(i)})$, using parameters $\boldsymbol{\theta}_d$. Analogously, let $\ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)}; \boldsymbol{\theta}_y)$ represent the loss of a classifier trained to predict the label $y^{(i)}$ from $\mathbf{f}(\mathbf{x}^{(i)})$, using parameters $\boldsymbol{\theta}_y$. The transformation \mathbf{f} can then be trained from two criteria: we want it to yield accurate predictions of the labels $y^{(i)}$, while making **inaccurate** predictions of the domains $d^{(i)}$. This can be formulated as a joint optimization problem,

$$\min_{\mathbf{f}, \boldsymbol{\theta}_y, \boldsymbol{\theta}_d} \sum_{i=1}^{N_\ell + N_u} \ell(\mathbf{f}(\mathbf{x}^{(i)}), d^{(i)}; \boldsymbol{\theta}_d) - \sum_{i=1}^{N_\ell} \ell(\mathbf{f}(\mathbf{x}^{(i)}), y^{(i)}; \boldsymbol{\theta}_y), \quad [4.40]$$

where N_ℓ is the number of labeled instances and N_u is the number of unlabeled instances, with the labeled instances appearing first in the dataset. The loss can be optimized by stochastic gradient descent, jointly training the parameters of the non-linear transformation \mathbf{f} and the parameters of the prediction models, $\boldsymbol{\theta}_d$ and $\boldsymbol{\theta}_y$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

4.5 *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: in practical applications, considerable attention may need to 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.

4.5.1 Sampling

In EM clustering, we maintain a distribution $q^{(i)}$ for the missing data related to each instance. In sampling-based algorithms, rather than maintaining a distribution over each latent variable, we draw random samples of the latent variables. If the sampling distribution is designed correctly, this procedure will eventually converge to drawing samples from the true posterior over the missing data, $p(\mathbf{z}^{(1:N_z)} \mid \mathbf{x}^{(1:N_x)})$. For example, in the case of clustering, the missing data $\mathbf{z}^{(1:N_z)}$ is the set of cluster memberships, $\mathbf{y}^{(1:N)}$, so we draw samples from the posterior distribution over clusterings of the data. If a single clustering is required, we can select the one with the highest conditional likelihood, $\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}} p(\mathbf{y}^{(1:N)} \mid \mathbf{x}^{(1:N)})$.

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 an MCMC algorithm in which each latent variable is sampled from its posterior distribution,

$$z^{(n)} \mid \mathbf{x}, \mathbf{z}^{(-n)} \sim p(z^{(n)} \mid \mathbf{x}, \mathbf{z}^{(-n)}), \quad [4.41]$$

where $\mathbf{z}^{(-n)}$ indicates $\{\mathbf{z} \setminus z^{(n)}\}$, the set of all latent variables except for $z^{(n)}$. By repeatedly sampling over all latent variables, we construct a Markov chain that converges to a set of samples from the posterior, $p(\mathbf{z}^{(1:N_z)} \mid \mathbf{x}^{(1:N_x)})$. In probabilistic clustering, the sampling distribution has the following form,

$$p(y^{(i)} \mid \mathbf{x}, \mathbf{y}^{(-i)}) = \frac{p(\mathbf{x}^{(i)} \mid y^{(i)}; \phi) \times p(y^{(i)}; \mu)}{\sum_{k=1}^K p(\mathbf{x}^{(i)} \mid k; \phi) \times p(k; \mu)} \quad [4.42]$$

$$\propto \text{Multinomial}(\mathbf{x}^{(i)}; \phi_{y^{(i)}}) \times \mu_{y^{(i)}}. \quad [4.43]$$

Note that in this case, the sampling distribution does not depend on the other instances $\mathbf{x}^{(-i)}, \mathbf{y}^{(-i)}$: given the parameters ϕ and μ , the posterior distribution over each $y^{(i)}$ can be computed from $\mathbf{x}^{(i)}$ alone.

(c) Jacob Eisenstein 2014-2017. Work in progress.

In sampling algorithms, there are several choices for how to deal with the parameters. One possibility is to sample them too. To do this, we must add them to the generative story, by introducing a prior distribution. For the multinomial and categorical parameters in the EM clustering model, the **Dirichlet distribution** is a typical choice, since it defines a probability on exactly the set of vectors that can be parameters: vectors that sum to one and include only non-negative numbers.¹⁰

To incorporate this prior, we augment the generative model with the following lines,

- For each cluster $k \in \{1, 2, \dots, K\}$,
 - Draw $\phi_k \sim \text{Dirichlet}(\alpha_\phi)$,
- Draw $\mu \sim \text{Dirichlet}(\alpha_\mu)$

where α is a parameter of the prior distribution, typically set to a constant vector $\alpha = [\alpha, \alpha, \dots, \alpha]$. When α is large, the Dirichlet distribution tends to generate vectors that are nearly uniform; when α is small, it tends to generate vectors that assign most of their probability mass to a few entries. Given prior distributions over ϕ and μ , we can now include them in Gibbs sampling, drawing values for these parameters from posterior distributions that are conditioned on the other variables in the model.

Unfortunately, sampling ϕ and μ can lead to slow convergence: the sampling distribution for these parameters is tightly constrained by the cluster memberships $\mathbf{y}^{(i)}$, which in turn is tightly constrained by the parameters. One alternative is to maintain ϕ and μ as parameters rather than latent variables. We can employ sampling in the E-step of the EM algorithm, constructing the distribution $\mathbf{q}^{(i)}$ from samples, and then updating the parameters using expected counts. This EM-MCMC hybrid is known as Monte Carlo Expectation Maximization (MCEM; Wei and Tanner, 1990), and is well-suited for cases in which it is difficult to compute $\mathbf{q}^{(i)}$ directly. Another possibility is to analytically integrate ϕ and μ out of the model. The cluster memberships $\mathbf{y}^{(i)}$ are the only remaining latent variable; we sample them from the compound distribution,

$$p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(1:N)}, \mathbf{y}^{(-i)}; \alpha_\phi, \alpha_\mu) = \int_{\phi, \mu} p(\phi, \mu \mid \mathbf{y}^{(-i)}, \mathbf{x}^{(1:N)}; \alpha_\phi, \alpha_\mu) p(\mathbf{y}^{(i)} \mid \mathbf{x}^{(1:N)}, \mathbf{y}^{(-i)}, \phi, \mu) d\phi d\mu. \quad [4.46]$$

¹⁰If $\sum_i^K \theta_i = 1$ and $\theta_i \geq 0$ for all i , then θ is said to be on the $K - 1$ **simplex**. A Dirichlet distribution with parameter $\alpha \in \mathbb{R}_+^K$ has support over the $K - 1$ simplex,

$$p_{\text{Dirichlet}}(\theta \mid \alpha) = \frac{1}{B(\alpha)} \prod_{i=1}^K \theta_i^{\alpha_i - 1} \quad [4.44]$$

$$B(\alpha) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i)}, \quad [4.45]$$

with $\Gamma(\cdot)$ indicating the gamma function, a generalization of the factorial function to non-negative reals.

(c) Jacob Eisenstein 2014-2017. Work in progress.

For multinomial and Dirichlet distributions, the sampling distribution can be computed in closed form. This is known as **collapsed Gibbs sampling**, because the parameters are “collapsed” out of the model.

MCMC algorithms are guaranteed to converge to the true posterior distribution over the latent variables, but there is no way to know how long this will take. In practice, the rate of convergence depends on initialization, just as expectation-maximization depends on initialization to avoid local optima. 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.

4.5.2 Spectral learning

Another approach to learning with latent variables is based on the **method of moments**, which makes it possible to avoid the problem of non-convex log-likelihood. Let us write $\bar{x}^{(i)}$ for the normalized vector of word counts in document i , so that $\bar{x}^{(i)} = \mathbf{x}^{(i)} / \sum_{j=1}^V x_j^{(i)}$. Then we can form a matrix of word-word co-occurrence counts,

$$\mathbf{C} = \sum_{i=1}^N \bar{\mathbf{x}}^{(i)} (\bar{\mathbf{x}}^{(i)})^\top. \quad [4.47]$$

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

$$E[\mathbf{C}] = \sum_{i=1}^N \sum_{k=1}^K \Pr(Z^{(i)} = k; \mu) \phi_k \phi_k^\top \quad [4.48]$$

$$= \sum_k^K N \mu_k \phi_k \phi_k^\top \quad [4.49]$$

$$= \Phi \text{Diag}(N\mu) \Phi^\top, \quad [4.50]$$

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,

$$\mathbf{C} = \Phi \text{Diag}(N\mu) \Phi^\top, \quad [4.51]$$

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**.

While moment-matching and eigendecomposition are similar in form, there is a key difference in the constraints on the solutions: in eigendecomposition, we require orthonormality, so that $\mathbf{Q} \mathbf{Q}^\top = \mathbb{I}$; in estimating the parameters of a mixture model, we require that

(c) Jacob Eisenstein 2014-2017. Work in progress.

μ and the columns of Φ are probability vectors. Thus, spectral learning algorithms must include a procedure for converting the solution into vectors that are non-negative and sum to one. 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 computing $p(z^{(i)} \mid \mathbf{x}^{(i)}; \phi, \mu) \propto p(\mathbf{x}^{(i)} \mid z^{(i)}; \phi)p(z^{(i)}; \mu)$. Spectral learning yields provably good solutions without regard to initialization, and can be quite fast in practice. However, it is more difficult to apply to a broad family of generative models than more generic techniques like EM and Gibbs Sampling. For more on applying spectral learning across a range of latent variable models, see Anandkumar et al. (2014).

Additional reading

There are a number of other learning paradigms that deviate from supervised learning.

- **Active learning:** the learner selects unlabeled instances and requests annotations (Settles, 2012).
- **Multiple instance learning:** labels are applied to bags of instances, with a positive label applied if at least one instance in the bag meets the criterion (Dietterich et al., 1997; Maron and Lozano-Pérez, 1998).
- **Constraint-driven learning:** supervision is provided in the form of explicit constraints on the learner (Chang et al., 2007; Ganchev et al., 2010).
- **Distant supervision:** noisy labels are generated from an external resource (Mintz et al., 2009).
- **Multitask learning:** the learner induces a representation that can be used to solve multiple classification tasks (Collobert et al., 2011b).
- **Transfer learning:** the learner must solve a classification task that differs from the labeled data (Pan and Yang, 2010).

Expectation maximization was introduced by Dempster et al. (1977), and is discussed in more detail by Murphy (2012). Like most machine learning treatments, Murphy focus on continuous observations and Gaussian likelihoods, rather than the discrete observations typically encountered in natural language processing. Murphy (2012) also includes an excellent chapter on MCMC; for an even more comprehensive treatment, see Robert and Casella (2013). For still more on Bayesian latent variable models, see Barber (2012), and for applications to natural language processing, see Cohen (2016). Surveys are available for semi-supervised learning (Zhu and Goldberg, 2009) and domain adaptation (Søgaard, 2013), although neither is recent enough to incorporate deep learning.

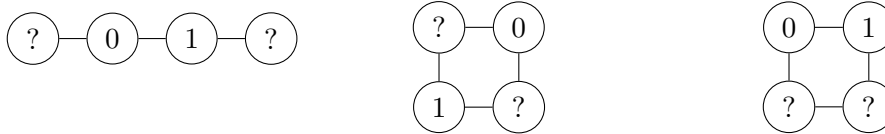
(c) Jacob Eisenstein 2014-2017. Work in progress.

Exercises

1. Derive the expectation maximization update for the parameter μ in the EM clustering model.
2. The expectation maximization lower bound \mathcal{J} is defined in Equation 4.10. Prove that the inverse $-\mathcal{J}$ is convex in q . You can use the following facts about convexity:
 - $f(x)$ is convex in x iff $\alpha f(x_1) + (1 - \alpha)f(x_2) \geq f(\alpha x_1 + (1 - \alpha)x_2)$ for all $\alpha \in [0, 1]$.
 - If $f(x)$ and $g(x)$ are both convex in x , then $f(x) + g(x)$ is also convex in x .
 - $\log(x + y) \leq \log x + \log y$.
3. Derive the E-step and M-step updates for the following generative model. You may assume that the labels $y^{(i)}$ are observed, but $z_m^{(i)}$ is not.
 - For each instance i ,
 - Draw label $y^{(i)} \sim \text{Categorical}(\mu)$
 - For each token $m \in \{1, 2, \dots, M^{(i)}\}$
 - * Draw $z_m^{(i)} \sim \text{Categorical}(\pi)$
 - * If $z_m^{(i)} = 0$, draw the current token from a label-specific distribution, $w_m^{(i)} \sim \phi_{y^{(i)}}$
 - * If $z_m^{(i)} = 1$, draw the current token from a document-specific distribution, $w_m^{(i)} \sim \nu^{(i)}$
4. Use expectation-maximization clustering to train a word-sense induction system, applied to the word *say*.
 - Import `nltk`, run `nltk.download()` and select `semcor`. Import `semcor` from `nltk.corpus`.
 - The command `semcor.tagged_sentences(tag='sense')` returns an iterator over sense-tagged sentences in the corpus. Each sentence can be viewed as an iterator over `tree` objects. For `tree` objects that are sense-annotated words, you can access the annotation as `tree.label()`, and the word itself with `tree.leaves()`. So `semcor.tagged_sentences(tag='sense')[0][2].label()` would return the sense annotation of the third word in the first sentence.
 - Extract all sentences containing the senses `say.v.01` and `say.v.02`.
 - Build bag-of-words vectors $x^{(i)}$, containing the counts of other words in those sentences, including all words that occur in at least two sentences.
 - Implement and run expectation-maximization clustering on the merged data.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Compute the frequency with which each cluster includes instances of `say.v.01` and `say.v.02`.
5. Using the iterative updates in Equations 4.34-4.36, compute the outcome of the label propagation algorithm for the following examples.



The value inside the node indicates the label, $y^{(i)} \in \{0, 1\}$, with $y^{(i)} = ?$ for unlabeled nodes. The presence of an edge between two nodes indicates $w_{i,j} = 1$, and the absence of an edge indicates $w_{i,j} = 0$. For the third example, you need only compute the first three iterations, and then you can guess at the solution in the limit.

In the remaining exercises, you will try out some approaches for semisupervised learning and domain adaptation. You will need datasets in multiple domains. You can obtain product reviews in multiple domains here: https://www.cs.jhu.edu/~mdredze/datasets/sentiment/processed_acl.tar.gz. Choose a source and target domain, e.g. dvds and books, and divide the data for the target domain into training and test sets of equal size.

6. First, quantify the cost of cross-domain transfer.
- Train a logistic regression classifier on the source domain training set, and evaluate it on the target domain test set.
 - Train a logistic regression classifier on the target domain training set, and evaluate it on the target domain test set. This is the “direct transfer” baseline.

Compute the difference in accuracy, which is a measure of the transfer loss across domains.

7. Next, apply the **label propagation** algorithm from § 4.3.2.

As a baseline, using only 5% of the target domain training set, train a classifier, and compute its accuracy on the target domain test set.

Next, apply label propagation:

- Compute the label matrix \mathbf{Q}_L for the labeled data (5% of the target domain training set), with each row equal to an indicator vector for the label (positive or negative).

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Iterate through the target domain instances, including both test and training data. At each instance i , compute all w_{ij} , using Equation 4.32, with $\alpha = 0.01$. Use these values to fill in column i of the transition matrix \mathbf{T} , setting all but the ten largest values to zero for each column i . Be sure to normalize the column so that the remaining values sum to one. You may need to use a sparse matrix for this to fit into memory.
- Apply the iterative updates from Equations 4.34-4.36 to compute the outcome of the label propagation algorithm for the unlabeled examples.

Select the test set instances from \mathbf{Q}_U , and compute the accuracy of this method. Compare with the supervised classifier trained only on the 5% sample of the target domain training set.

- Using only 5% of the target domain training data (and all of the source domain training data), implement one of the supervised domain adaptation baselines in § 4.4.1. See if this improves on the “direct transfer” baseline from the previous problem
- Implement `EasyAdapt` (§ 4.4.1), again using 5% of the target domain training data and all of the source domain data.
- Now try unsupervised domain adaptation, using the “linear projection” method described in § 4.4.2. Specifically:
 - Identify 500 pivot features as the words with the highest frequency in the (complete) training data for the source and target domains. Specifically, let x_i^d be the count of the word i in domain d : choose the 500 words with the largest values of $\min(x_i^{\text{source}}, x_i^{\text{target}})$.
 - Train a classifier to predict each pivot feature from the remaining words in the document.
 - Arrange the features of these classifiers into a matrix Φ , and perform truncated singular value decomposition, with $k = 20$
 - Train a classifier from the source domain data, using the combined features $\mathbf{x}^{(i)} \oplus \mathbf{U}^\top \mathbf{x}^{(i)}$ — these include the original bag-of-words features, plus the projected features.
 - Apply this classifier to the target domain test set, and compute the accuracy.

Part II

Sequences and trees

Chapter 5

Language models

In probabilistic classification, we are interested in computing the probability of a label, conditioned on the text. Let us now consider something like the inverse problem: computing the probability of text itself. Specifically, we will consider models that assign probability to a sequence of word tokens,¹ $p(w_1, w_2, \dots, w_M)$, with $w_m \in \mathcal{V}$. The set \mathcal{V} is a discrete vocabulary,

$$\mathcal{V} = \{\text{aardvark}, \text{abacus}, \dots, \text{zither}\}. \quad [5.1]$$

Why would we want to compute the probability of a word sequence? In many applications, our goal is to produce word sequences as output:

- In **machine translation**, we convert from text in a source language to text in a target language.
- In **speech recognition**, we convert from audio signal to text.
- In **summarization**, we convert from long texts into short texts.
- In **dialogue systems**, we convert from the user’s input (and perhaps an external knowledge base) into a text response.

In each of these cases, a key subcomponent is to compute the probability of the output text. By choosing high-probability output, we hope to generate texts that are more **fluent**. For example, suppose we want to translate a sentence from Spanish to English.

¹The linguistic term “word” does not cover everything we might want to model, such as names, numbers, and emoticons. Instead, we prefer the term **token**, which refers to anything that can appear in a sequence of linguistic data. **Tokenizers** are programs for segmenting strings of characters or bytes into tokens. In standard written English, tokenization is relatively straightforward, and can be performed using a regular expression. But in languages like Chinese, tokens are not usually separated by spaces, so tokenization can be considerably more challenging. For more on tokenization algorithms, see Manning et al. (2008), chapter 2.

(5.1) *El cafe negro me gusta mucho.*

A literal word-for-word translation (sometimes called a **gloss**) is,

(5.2) *The coffee black me pleases much.*

A good language model of English will tell us that the probability of this translation is low, in comparison with more grammatical alternatives, such as,

$$p(\textit{The coffee black me pleases much}) < p(\textit{I love dark coffee}). \quad [5.2]$$

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 (Weaver, 1955):

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 observation motivates a generative model (like Naïve Bayes):

- The English sentence $w^{(e)}$ is generated from a **language model**, $p_e(w^{(e)})$.
- The Spanish sentence $w^{(s)}$ is then generated from a **translation model**, $p_{s|e}(w^{(s)} | w^{(e)})$.

Given these two distributions, we can then perform translation by Bayes rule:

$$p_{e|s}(w^{(e)} | w^{(s)}) \propto p_{e,s}(w^{(e)}, w^{(s)}) \quad [5.3]$$

$$= p_{s|e}(w^{(s)} | w^{(e)}) \times p_e(w^{(e)}). \quad [5.4]$$

This is sometimes called the **noisy channel model**, because it envisions English text turning into Spanish by passing through a noisy channel, $p_{s|e}$. What is the advantage of modeling translation this way, as opposed to modeling $p_{e|s}$ directly? The crucial point is that the two distributions $p_{s|e}$ (the translation model) and p_e (the language model) can be estimated from separate data. The translation model requires **bitext** — examples of correct translations. But the language model requires only text in English. Such monolingual data is much more widely available, which means that the fluency of the output translation can be improved simply by scraping more webpages. Furthermore, once estimated, the language model p_e can be reused in any application that involves generating English text, from summarization to speech recognition.

(c) Jacob Eisenstein 2014-2017. Work in progress.

5.1 N-gram language models

How can we estimate the probability of a sequence of word tokens? The simplest idea would be to apply a **relative frequency estimator**. For example, consider the quote, attributed to Picasso, “*computers are useless, they can only give you answers.*” We can estimate the probability of this sentence as follows:

$$\begin{aligned} p(\text{Computers are useless, they can only give you answers}) \\ = \frac{\text{count}(\text{Computers are useless, they can only give you answers})}{\text{count}(\text{all sentences ever spoken})} \end{aligned} \quad [5.5]$$

This estimator is **unbiased**: in the theoretical limit of infinite data, the estimate will be correct. But in practice, we are asking for accurate counts over an infinite number of events, since sequences of words can be arbitrarily long. Even if we set an aggressive upper bound of, say, $n = 20$ tokens in the sequence, the number of possible sequences is V^{20} . A small vocabulary for English would have $V = 10^4$, so we would have 10^{80} possible sequences. Clearly, this estimator is very data-hungry, and suffers from high variance: even grammatical sentences will have probability zero if they happen not to have occurred in the training data.² We therefore need to introduce bias to have a chance of making reliable estimates from finite training data. The language models that follow in this chapter introduce bias in various ways.

We begin with n -gram language models, which compute the probability of a sequence as the product of probabilities of subsequences. The probability of a sequence $p(\mathbf{w}) = p(w_1, w_2, \dots, w_M)$ can be refactored using the chain rule:

$$p(\mathbf{w}) = p(w_1, w_2, \dots, w_M) \quad [5.6]$$

$$= p(w_1) \times p(w_2 \mid w_1) \times p(w_3 \mid w_2, w_1) \times \dots \times p(w_M \mid w_{M-1}, \dots, w_1) \quad [5.7]$$

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: given the context *Computers are*, we want to compute a probability over the next token. The relative frequency estimate of the probability of the word *useless* in this context is,

$$\begin{aligned} p(\text{useless} \mid \text{computers are}) &= \frac{\text{count}(\text{computers are useless})}{\sum_{x \in \mathcal{V}} \text{count}(\text{computers are } x)} \\ &= \frac{\text{count}(\text{computers are useless})}{\text{count}(\text{computers are})}. \end{aligned}$$

²Chomsky has famously argued that this is evidence against the very concept of probabilistic language models: no such model could distinguish the grammatical sentence *colorless green ideas sleep furiously* from the ungrammatical permutation *furiously sleep ideas green colorless*. Indeed, even the bigrams in these two examples are unlikely to occur — at least, not in texts written before Chomsky proposed this example.

Note that we haven't made any approximations yet, and we could have just as well applied the chain rule in reverse order,

$$p(\mathbf{w}) = p(w_M) \times p(w_{M-1} \mid w_M) \times \dots \times p(w_1 \mid w_2, \dots, w_M), \quad [5.8]$$

or in any other order. But this means that we also haven't really improved anything either: to compute the conditional probability $p(w_M \mid w_{M-1}, w_{M-2}, \dots, w_1)$, we need to model V^{M-1} contexts. We cannot estimate such a distribution from any reasonable finite sample.

To solve this problem, n -gram models make a crucial simplifying approximation: condition on only the past $n - 1$ words.

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

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

$$p(w_1, \dots, w_M) \approx \prod_m^M p(w_m \mid w_{m-1}, \dots, w_{m-n+1}) \quad [5.10]$$

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

$$p(I \text{ like black coffee}) = p(I \mid \diamond) \times p(\text{like} \mid I) \times p(\text{black} \mid \text{like}) \times p(\text{coffee} \mid \text{black}) \times p(\blacklozenge \mid \text{coffee}). \quad [5.11]$$

In this model, we have to estimate and store the probability of only V^n events, which is exponential in the order of the n -gram, and not V^M , which is exponential in the length of the sentence. The n -gram probabilities can be computed by relative frequency estimation,

$$\Pr(W_m = i \mid W_{m-1} = j, W_{m-2} = k) = \frac{\text{count}(i, j, k)}{\sum_{i'} \text{count}(i', j, k)} = \frac{\text{count}(i, j, k)}{\text{count}(j, k)} \quad [5.12]$$

A key design question is how to set the hyperparameter n , which controls the size of the context used in each conditional probability. If this is misspecified, the language model will sacrifice accuracy. Let's consider the potential problems concretely.

When n is too small. Consider the following sentences:

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

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

(c) Jacob Eisenstein 2014-2017. Work in progress.

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.

When n is too big. In this case, we cannot 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. A small n -gram size introduces high bias with respect to the true distribution, and a large n -gram size introduces high variance due to the huge number of possible events. But in reality the situation is even worse, because 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**.

5.2 Smoothing and discounting

Limited data is a persistent problem in estimating language models. In § 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).

5.2.1 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 imaginary “pseudo” counts. The same idea can be applied to n -gram language models, as shown here in the bigram case,

$$P_{\text{smooth}}(w_m \mid w_{m-1}) = \frac{\text{count}(w_{m-1}, w_m) + \alpha}{\sum_{w' \in \mathcal{V}} \text{count}(w_{m-1}, w') + V\alpha}. \quad [5.13]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

This basic framework is called **Lidstone smoothing**, but special cases have other names:

- **Laplace smoothing** corresponds to the case $\alpha = 1$.
- **Jeffreys-Perks law** corresponds to the case $\alpha = 0.5$. Manning and Schütze (1999) offer more insight on the justifications for this setting.

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

$$c_i^* = (c_i + \alpha) \frac{M}{M + V\alpha}, \quad [5.14]$$

where c_i is the count of event i , c_i^* is the effective count, and $M = \sum_i^V c_i$ is the total number of terms in the dataset (w_1, w_2, \dots, w_M) . This term ensures that $\sum_i^V c_i^* = \sum_i^V c_i = M$. The **discount** for each n-gram is then computed as,

$$d_i = \frac{c_i^*}{c_i} = \frac{(c_i + \alpha)}{c_i} \frac{M}{(M + V\alpha)}.$$

5.2.2 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, suppose we set an absolute discount $d = 0.1$ in a bigram model, and then redistribute this probability mass equally over the unseen words. The resulting probabilities are shown in Table 5.1.

Discounting reserves some probability mass from the observed data, and we need not redistribute this probability mass equally. Instead, we can **backoff** 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**. In this smoothing model, bigram probabilities are computed as,

$$c^*(i, j) = c(i, j) - d \quad [5.15]$$

$$P_{\text{Katz}}(i | j) = \begin{cases} \frac{c^*(i, j)}{c(j)} & \text{if } c(i, j) > 0 \\ \alpha(j) \times \frac{P_{\text{unigram}}(i)}{\sum_{i': c(i', j)=0} P_{\text{unigram}}(i')} & \text{if } c(i, j) = 0. \end{cases} \quad [5.16]$$

The term $\alpha(j)$ indicates the amount of probability mass that has been discounted for context j . This probability mass is then divided across all the unseen events, $\{i' : c(i', j) =$

			Lidstone smoothing, $\alpha = 0.1$		Discounting, $d = 0.1$	
	counts	unsmoothed probability	effective counts	smoothed probability	effective counts	smoothed probability
<i>impropriety</i>	8	0.4	7.826	0.391	7.9	0.395
<i>offense</i>	5	0.25	4.928	0.246	4.9	0.245
<i>damage</i>	4	0.2	3.961	0.198	3.9	0.195
<i>deficiencies</i>	2	0.1	2.029	0.101	1.9	0.095
<i>outbreak</i>	1	0.05	1.063	0.053	0.9	0.045
<i>infirmity</i>	0	0	0.097	0.005	0.25	0.013
<i>cephalopods</i>	0	0	0.097	0.005	0.25	0.013

Table 5.1: Example of Lidstone smoothing and absolute discounting in a bigram language model, for the context (*alleged*, -), for a toy corpus with a total of twenty counts over the seven words shown. Note that discounting decreases the probability for all but the unseen words, while Lidstone smoothing increases the effective counts and probabilities for *deficiencies* and *outbreak*.

0}, proportional to the unigram probability of each word i' . The discount parameter d can be optimized to maximize performance (typically held-out log-likelihood) on a development set.

5.2.3 *Interpolation

Backoff is one way to combine n -gram models across various values of n . An alternative approach 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. For example, for an interpolated trigram model,

$$\begin{aligned}
 P_{\text{Interpolation}}(i \mid j, k) &= \lambda_3 p_3^*(i \mid j, k) \\
 &\quad + \lambda_2 p_2^*(i \mid j) \\
 &\quad + \lambda_1 p_1^*(i).
 \end{aligned}$$

In this equation, p_n^* is the unsmoothed empirical probability given by an n -gram language model, and λ_n is the weight assigned to this model. To ensure that the interpolated $p(w)$ is still a valid probability distribution, we must obey the constraint, $\sum_n \lambda_n = 1$. But how to find the specific values of λ ?

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. What's missing in this case? We can think of each word w_m as drawn from an n -gram of unknown size, $z_m \in \{1 \dots n_{\max}\}$. This z_m is the missing data that we are looking for. Therefore, the application of EM to this problem involves the following **generative process**:

- For each token $m \in \{1, 2, \dots, M\}$:
 - draw $z_m \sim \text{Categorical}(\lambda)$,
 - draw $w_m \sim p_{z_m}^*(w_m \mid w_{m-1}, \dots, w_{m-z_m})$.

If the missing data $\{Z_m\}$ were known, then we could estimate λ from relative frequency estimation,

$$\lambda_z = \frac{\text{count}(Z_m = z)}{M} \quad [5.17]$$

$$\propto \sum_{m=1}^M \delta(Z_m = z). \quad [5.18]$$

But since we do not know the values of the latent variables Z_m , we impute a distribution q_m in the E-step, which represents the degree of belief that word token w_m was generated from a n -gram of order z_m ,

$$q_m(z) \triangleq \Pr(Z_m = z \mid \mathbf{w}_{1:m}; \lambda) \quad [5.19]$$

$$= \frac{p(w_m \mid \mathbf{w}_{1:m-1}, Z_m = z) \times p(z)}{\sum_{z'} p(w_m \mid \mathbf{w}_{1:m-1}, Z_m = z') \times p(z')} \quad [5.20]$$

$$\propto p_z^*(w_m \mid \mathbf{w}_{1:m-1}) \times \lambda_z. \quad [5.21]$$

In the M-step, we can compute λ by summing the expected counts under q ,

$$\lambda_z \propto \sum_{m=1}^M q_m(z). \quad [5.22]$$

By iterating between updates to q and λ , we will ultimately converge at a solution. The complete algorithm is shown in Algorithm 9.

5.2.4 *Kneser-Ney smoothing

Kneser-Ney smoothing is based on absolute discounting, but it redistributes the resulting probability mass in a different way from Katz backoff. Empirical evidence points to Kneser-Ney smoothing as the state-of-art for n -gram language modeling (Goodman, 2001). To motivate Kneser-Ney smoothing, consider the example: *I recently visited ..* Which of the following is more likely?

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 9 Expectation-maximization for interpolated language modeling

```

1: procedure ESTIMATE INTERPOLATED  $n$ -GRAM ( $\mathbf{w}_{1:M}, \{\mathbf{p}_n^*\}_{n \in 1:n_{\max}}$ )
2:   for  $z \in \{1, 2, \dots, n_{\max}\}$  do ▷ Initialization
3:      $\lambda_z \leftarrow \frac{1}{n_{\max}}$ 
4:   repeat
5:     for  $m \in \{1, 2, \dots, M\}$  do ▷ E-step
6:       for  $z \in \{1, 2, \dots, n_{\max}\}$  do
7:          $q_m(z) \leftarrow \mathbf{p}_z^*(w_m \mid \mathbf{w}_{1:m-}) \times \lambda_z$ 
8:        $\mathbf{q}_m \leftarrow \text{Normalize}(\mathbf{q}_m)$ 
9:       for  $z \in \{1, 2, \dots, n_{\max}\}$  do ▷ M-step
10:         $\lambda_z \leftarrow \frac{1}{M} \sum_{m=1}^M q_m(z)$ 
11:   until tired
12:   return  $\lambda$ 

```

- *Francisco*
- *Duluth*

Now suppose that both bigrams *visited Duluth* and *visited Francisco* are unobserved in our training data, and furthermore, that the unigram probability $\mathbf{p}_1^*(\textit{Francisco})$ is greater than $\mathbf{p}^*(\textit{Duluth})$. Nonetheless we would still guess that $\mathbf{p}(\textit{visited Duluth}) > \mathbf{p}(\textit{visited Francisco})$, because *Duluth* is a more **versatile** word: it an occur in many contexts, while *Francisco* usually occurs in a single context, following the word *San*. This notion of versatility is the key to Kneser-Ney smoothing.

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

$$\mathbf{p}_{KN}(w \mid u) = \begin{cases} \frac{\text{count}(w, u) - d}{\text{count}(u)}, & \text{count}(w, u) > 0 \\ \alpha(u) \times \mathbf{p}_{\text{continuation}}(w), & \text{otherwise} \end{cases} \quad [5.23]$$

$$\mathbf{p}_{\text{continuation}}(w) = \frac{|u : \text{count}(w, u) > 0|}{\sum_{w' \in \mathcal{V}} |u' : \text{count}(w', u') > 0|}. \quad [5.24]$$

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 u is $d \times |w : \text{count}(w, u) > 0|$, and we divide this probability mass equally among the unseen n -grams,

$$\alpha(u) = |w : \text{count}(w, u) > 0| \times \frac{d}{\text{count}(u)}. \quad [5.25]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

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 u in which w appears, and in the denominator, we normalize by summing the same quantity over all words w' .

The idea of modeling versatility by counting contexts may seem heuristic, but there is an elegant theoretical justification from Bayesian nonparametrics (Teh, 2006). Kneser-Ney smoothing on n -grams was the dominant language modeling technique — widely used in speech recognition and machine translation — before the arrival of neural language models.

5.3 Recurrent neural network language models

Until this decade, n -grams were the dominant language modeling approach. But in a few years, they have been almost completely supplanted by a new family of language models based on **neural networks**. These models do not make the n -gram assumption of restricted context; indeed, they can incorporate arbitrarily distant contextual information, while remaining computationally and statistically tractable.

The first insight is to treat word prediction as a **discriminative** learning task: rather than directly estimating the distribution $p(w \mid u)$ from (smoothed) relative frequencies, we now treat language modeling as a machine learning problem, and estimate parameters that maximize the log conditional probability of a corpus.³

The second insight is to reparametrize the probability distribution $p(w \mid u)$ as a function of two dense K -dimensional numerical vectors, $\beta_w \in \mathbb{R}^K$, and $v_u \in \mathbb{R}^K$,

$$p(w \mid u) = \frac{\exp(\beta_w \cdot v_u)}{\sum_{w' \in \mathcal{V}} \exp(\beta_{w'} \cdot v_u)}, \quad [5.26]$$

where $\beta_w \cdot v_u$ represents a dot product. Note that the denominator ensures that it is a properly normalized probability distribution. In the neural networks literature, this function is sometimes known as a **softmax** layer, written

$$(\text{SoftMax}(\mathbf{a}))_i = \frac{\exp(a_i)}{\sum_j \exp(a_j)}, \quad [5.27]$$

where \mathbf{a} is a vector of scores and $\text{SoftMax}(\mathbf{a})_i$ is a normalized distribution.⁴

³This idea is not in itself new; for example, Rosenfeld (1996) applies logistic regression to language modeling, and Roark et al. (2007) apply perceptrons and conditional random fields (§ 6.5.3).

⁴The logistic regression classifier can be viewed as an application of the softmax transformation to the vector constructed by computing the inner products of weights and features for all possible labels.

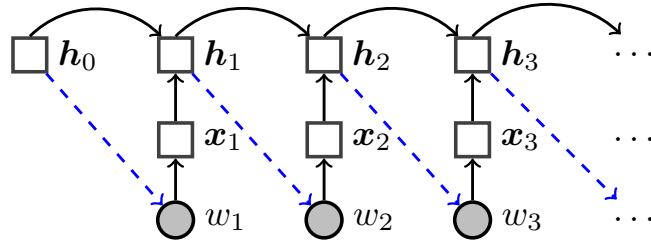


Figure 5.1: The recurrent neural network language model, viewed as an “unrolled” computation graph. Solid lines indicate direct computation, dotted blue lines indicate probabilistic dependencies, circles indicate random variables, and squares indicate computation nodes.

The word vectors β_w are parameters of the model, and are estimated directly. As we will see in chapter 13, these vectors carry useful information about word meaning, and semantically similar words tend to have highly correlated vectors.

The context vectors v_u can be computed in various ways, depending on the model. Here we will consider a relatively simple — but effective — neural language model, the **recurrent neural network** (RNN; Mikolov et al., 2010). The basic idea is to recurrently update the context vectors as we move through the sequence. Let us write h_m for the contextual information at position m in the sequence. RNNs employ the following recurrence:

$$x_m \triangleq \phi_{w_m} \quad [5.28]$$

$$h_m = g(\Theta h_{m-1} + x_m) \quad [5.29]$$

$$p(w_{m+1} \mid w_1, w_2, \dots, w_m) = \frac{\exp(\beta_{w_{m+1}} \cdot h_m)}{\sum_{w' \in \mathcal{V}} \exp(\beta_{w'} \cdot h_m)}, \quad [5.30]$$

where ϕ is a matrix of **input word embeddings**, and x_m denotes the embedding for word w_m . The function g is an element-wise nonlinear **activation function**. Typical choices are:

- $\tanh(x)$, the hyperbolic tangent;
- $\sigma(x)$, the **sigmoid function** $\frac{1}{1+\exp(-x)}$;
- $(x)_+$, the **rectified linear unit**, $(x)_+ = \max(x, 0)$, also called **ReLU**.

These activation functions are shown in Figure 5.2. The sigmoid and tanh functions “squash” their inputs into a fixed range: $[0, 1]$ for the sigmoid, $[-1, 1]$ for tanh. This makes it possible to chain together many instances of these functions without numerical instability.

(c) Jacob Eisenstein 2014-2017. Work in progress.

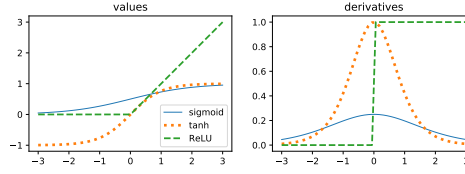


Figure 5.2: Nonlinear activation functions for neural networks

A key point about the RNN language model is that although each w_m depends only on the context vector \mathbf{h}_{m-1} , this vector is in turn influenced by **all** previous tokens, w_1, w_2, \dots, w_{m-1} , through the recurrence operation: w_1 affects \mathbf{h}_1 , which affects \mathbf{h}_2 , and so on, until the information is propagated all the way to \mathbf{h}_{m-1} , and then on to w_m (see Figure 5.1). This is an important distinction from n -gram language models, where any information outside the n -word window is ignored. Thus, in principle, the RNN language model can handle long-range dependencies, such as number agreement over long spans of text — although it would be difficult to know where exactly in the vector \mathbf{h}_m this information is represented. The main limitation is that information is attenuated by repeated application of the nonlinearity g . **Long short-term memories** (LSTMs), described below, are a variant of RNNs that address this issue, using memory cells to propagate information through the sequence without applying non-linearities (Hochreiter and Schmidhuber, 1997).

The denominator in Equation 5.30 is a computational bottleneck, because it involves a sum over the entire vocabulary. One solution is to use a **hierarchical softmax** function, which computes the sum more efficiently by organizing the vocabulary into a tree (Mikolov et al., 2011). Another strategy is to optimize an alternative metric, such as **noise-contrastive estimation** (Gutmann and Hyvärinen, 2012), which learns by distinguishing observed instances from artificial instances generated from a noise distribution (Mnih and Teh, 2012).

5.3.1 Estimation by backpropagation

The recurrent neural network language model has the following parameters:

- $\phi_i \in \mathbb{R}^K$, the “input” word vectors (these are sometimes called **word embeddings**, since each word is embedded in a K -dimensional space);
- $\beta_i \in \mathbb{R}^K$, the “output” word vectors;
- $\Theta \in \mathbb{R}^{K \times K}$, the recurrence operator.

Each of these parameters must be estimated. We do this by formulating an objective function over the training corpus, $L(w)$, and then employ **backpropagation** to incrementally update the parameters after encountering each training example. Backpropagation is a term from the neural network literature, which means that we use the chain rule of

differentiation to obtain gradients on each parameter. After obtaining these gradients, we can apply an online learning algorithm such as stochastic gradient descent (see § 1.5.2).

For example, suppose we want to obtain the gradient of the log-likelihood with respect to a single row of the recurrence operator, θ_k . Let us first define the total objective as a sum over local error functions e_m ,

$$\ell(\mathbf{w}) = \sum_{m=1}^M e_m(\mathbf{h}_{m-1}) \quad [5.31]$$

$$e_m(\mathbf{h}_{m-1}) \triangleq -\log \mathbf{p}(w_m \mid w_1, w_2, \dots, w_{m-1}) \quad [5.32]$$

$$= -\beta_{w_m} \cdot \mathbf{h}_{m-1} + \log \sum_{w' \in \mathcal{V}} \exp(\beta_{w'} \cdot \mathbf{h}_{m-1}). \quad [5.33]$$

We can now differentiate the objective with respect to θ_k :

$$\frac{\partial}{\partial \theta_k} \ell(\mathbf{w}) = \sum_m \frac{\partial e_m(\mathbf{h}_{m-1})}{\partial \theta_k} \quad [5.34]$$

$$= \sum_{m=1}^M (\nabla_{\mathbf{h}_{m-1}} e_m) \frac{\partial}{\partial \theta_k} \mathbf{h}_{m-1}. \quad [5.35]$$

In the first line, we simply distribute the derivative across the sum. In the second line, we apply the chain rule of calculus. The term $\nabla_{\mathbf{h}_{m-1}} e_m$ refers to the gradient of the error e_m evaluated at \mathbf{h}_{m-1} , and is equal to,

$$\nabla_{\mathbf{h}_{m-1}} e_m = -\beta_{w_m} + B \text{SoftMax}(B\mathbf{h}_{m-1}), \quad [5.36]$$

where B is a matrix with all word output embeddings stacked vertically, $B = (\beta_1^\top, \beta_2^\top, \dots, \beta_V^\top)$.

Next we compute the derivative of \mathbf{h}_{m-1} , first noting that within the vector \mathbf{h}_{m-1} , only the element $h_{m-1,k}$ depends on θ_k .

$$\frac{\partial}{\partial \theta_k} \mathbf{h}_{m-1} = \frac{\partial}{\partial \theta_k} h_{m-1,k} \quad [5.37]$$

$$= \frac{\partial}{\partial \theta_k} g(\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}) \quad [5.38]$$

$$= (\nabla_{\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}} g) \times \frac{\partial}{\partial \theta_k} (\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}) \quad [5.39]$$

$$= (\nabla_{\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}} g) \times (\mathbf{h}_{m-2} + \theta_k \odot \frac{\partial}{\partial \theta_k} \mathbf{h}_{m-2}), \quad [5.40]$$

where \odot is an elementwise (Hadamard) vector product, and $(\nabla_{\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}} g)$ is the elementwise gradient of the non-linear activation function for \mathbf{h}_{m-1} evaluated at the scalar

$\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}$. For example, if g is the elementwise hyperbolic tangent, then its gradient is,

$$(\nabla_{\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k}} g) = (1 - \tanh^2(\theta_k \cdot \mathbf{h}_{m-2} + x_{m-1,k})). \quad [5.41]$$

The application of backpropagation to sequence models such as recurrent neural networks is known as **backpropagation through time**. A key point is that the derivative $\frac{\partial \mathbf{h}_{m-1}}{\partial \theta_k}$ depends recurrently on $\frac{\partial \mathbf{h}_{m-2}}{\partial \theta_k}$, and on all $\frac{\partial \mathbf{h}_n}{\partial \theta_k}$ for $n < m$. Furthermore, we will need to compute $\frac{\partial \mathbf{h}_{m-2}}{\partial \theta_k}$ **again**, to account for the error term $e_{m-1}(\mathbf{h}_{m-2})$. To avoid redoing work, it is best to cache such derivatives, so that they can be reused during backpropagation.

Backpropagation is implemented by neural network toolkits such as TensorFlow (Abadi et al., 2016), Torch (Collobert et al., 2011a), and DyNet (Neubig et al., 2017a). In these toolkits, the user defines a **computation graph** representing the neural network structure, which culminates in a scalar loss function. The toolkit then automatically computes the gradient of the loss function with respect to all model parameters, by applying the chain rule of differentiation across the computation graph. Unlike the classification objectives considered in § 1.2, neural network objectives are usually non-convex function of the parameters, so there is no learning procedure that is guaranteed to converge to the global optimum. Nonetheless, gradient-based optimization often yields parameter estimates that are very effective in practice.

5.3.2 Hyperparameters

The RNN language model has several hyperparameters that must be tuned to ensure good performance. The model capacity is controlled by the size of the word and context vectors K , which play a role that is somewhat analogous to the size of the n -gram context. For datasets that are large with respect to the vocabulary (i.e., there is a large token-to-type ratio), we can afford to estimate a model with a large K , which enables more subtle distinctions between words and contexts. When the dataset is relatively small, then K must be smaller too. However, this general advice has not yet been formalized into any concrete formula for choosing K , and trial-and-error is still necessary. Overfitting can also be prevented by **dropout**, which involves randomly setting some elements of the computation to zero (Srivastava et al., 2014), forcing the learner not to rely too much on any particular dimension of the word or context vectors. (The dropout rate must also be tuned by the user.) Other design decisions include: the nature of the nonlinear activation function g , the size of the vocabulary, and the parametrization of the learning algorithm, such as the learning rate.

(c) Jacob Eisenstein 2014-2017. Work in progress.

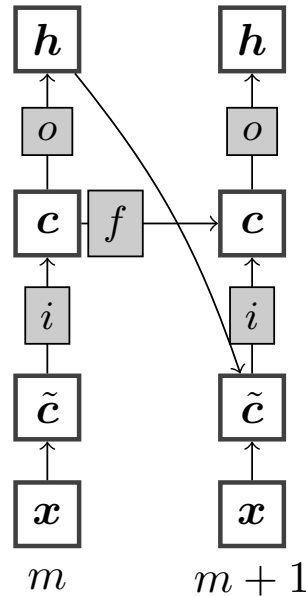


Figure 5.3: The long short-term memory (LSTM) architecture. In an LSTM language model, each h_m would be used to predict the next word w_{m+1} . [todo: needs work]

5.3.3 Alternative neural language models

A well known problem with RNNs is that backpropagation across long chains tends to lead to “vanishing” or “exploding” gradients (Bengio et al., 1994). For example, the input embedding of word w_1 affects the likelihood of a distant word such as w_{29} , but this impact may be attenuated by backpropagation through the intervening time steps. One solution is to rescale the gradients, or to clip them at some maximum value (Pascanu et al., 2013). An alternative is to change the model architecture itself.

A popular variant of RNNs, which is more robust to these problems, is the **long short-term memory (LSTM)** (Hochreiter and Schmidhuber, 1997; Sundermeyer et al., 2012). This model augments the hidden state h_m with a “memory cell” c_m . The value of the memory cell at each time m is a linear interpolation between two quantities: its previous value c_{m-1} , and an “update” \tilde{c}_m , which is computed from the current input x_m and the previous hidden state h_{m-1} . The next state h_m is then computed from the memory cell. Because the memory cell is never passed through the non-linear function g , it is possible for information to propagate through the network over long distances.

The interpolation weights are controlled by a set of gates, which are themselves functions of the input and previous hidden state. The gates are computed from sigmoid activations, ensuring that their values will be in the range $[0, 1]$. They can therefore be viewed

as soft, differentiable logic gates. The LSTM architecture is shown in Figure 5.3, and the complete update equations are:

$$\mathbf{f}_m = \sigma(\Theta^{(h \rightarrow f)} \cdot \mathbf{h}_{m-1} + \Theta^{(x \rightarrow f)} \cdot \mathbf{x}_m) \quad \text{forget gate} \quad [5.42]$$

$$\mathbf{i}_m = \sigma(\Theta^{(h \rightarrow i)} \cdot \mathbf{h}_{m-1} + \Theta^{(x \rightarrow i)} \cdot \mathbf{x}_m) \quad \text{input gate} \quad [5.43]$$

$$\tilde{\mathbf{c}}_m = \tanh(\Theta^{(h \rightarrow c)} \cdot \mathbf{h}_{m-1} + \Theta^{(w \rightarrow c)} \cdot \mathbf{x}_m) \quad \text{update candidate} \quad [5.44]$$

$$\mathbf{c}_m = \mathbf{f}_m \odot \mathbf{c}_{m-1} + \mathbf{i}_m \odot \tilde{\mathbf{c}}_m \quad \text{memory cell update} \quad [5.45]$$

$$\mathbf{o}_m = \sigma(\Theta^{(h \rightarrow o)} \cdot \mathbf{h}_{m-1} + \Theta^{(x \rightarrow o)} \cdot \mathbf{x}_m) \quad \text{output gate} \quad [5.46]$$

$$\mathbf{h}_m = \mathbf{o}_m \odot \mathbf{c}_m \quad \text{output.} \quad [5.47]$$

As above, \odot refers to an elementwise (Hadamard) product. The LSTM model has been shown to outperform standard recurrent neural networks across a wide range of problems (it was first used for language modeling by Sundermeyer et al. (2012)), and is now widely used for sequence modeling tasks. There are several LSTM variants, of which the Gated Recurrent Unit (Cho et al., 2014) is presently one of the more well known. Many software packages implement a variety of RNN architectures, so choosing between them is simple from a user’s perspective. Jozefowicz et al. (2015) provide an empirical comparison of various modeling choices circa 2015. Notable earlier non-recurrent architectures include the neural probabilistic language model (Bengio et al., 2003) and the log-bilinear language model (Mnih and Hinton, 2007). Much more detail on these models can be found in the text by Goodfellow et al. (2016).

5.4 Evaluating language models

Because language models are typically components of larger systems — language modeling is not usually 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.

5.4.1 Held-out likelihood

The goal of probabilistic language models is to accurately measure the probability of sequences of word tokens. Therefore, an intrinsic evaluation metric is the likelihood that the

(c) Jacob Eisenstein 2014-2017. Work in progress.

language model assigns to **held-out data**, which is not used during training. Specifically, we compute,

$$\ell(\mathbf{w}) = \sum_{m=1}^M \log p(w_m \mid w_{m-1}, \dots), \quad [5.48]$$

treating the entire held-out corpus as a single stream of tokens.

Typically, unknown words 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 $\langle \text{UNK} \rangle$. Other strategies for dealing with out-of-vocabulary terms are discussed in § 5.5.

5.4.2 Perplexity

Held-out likelihood is usually presented as **perplexity**, which is a deterministic transformation of the log-likelihood into an information-theoretic quantity,

$$\text{Perplex}(\mathbf{w}) = 2^{-\frac{\ell(\mathbf{w})}{M}}, \quad [5.49]$$

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

Lower perplexities correspond to higher likelihoods, so lower scores are better on this metric. (How to remember: lower perplexity is better, because your language model is less perplexed.) To understand perplexity, here are some special cases:

- In the limit of a perfect language model, probability 1 is assigned to the held-out corpus, with $\text{Perplex}(\mathbf{w}) = 2^{-\frac{1}{M} \log_2 1} = 2^0 = 1$.
- In the opposite limit, probability zero is assigned to the held-out corpus, which corresponds to an infinite perplexity, $\text{Perplex}(\mathbf{w}) = 2^{-\frac{1}{M} \log_2 0} = 2^\infty = \infty$.
- Assume a uniform, unigram model in which $p(w_i) = \frac{1}{V}$ for all words in the vocabulary. Then,

$$\begin{aligned} \log_2(\mathbf{w}) &= \sum_{m=1}^M \log_2 \frac{1}{V} = - \sum_{m=1}^M \log_2 V = -M \log_2 V \\ \text{Perplex}(\mathbf{w}) &= 2^{\frac{1}{M} M \log_2 V} \\ &= 2^{\log_2 V} \\ &= V. \end{aligned}$$

This is the “worst reasonable case” scenario, since you could build such a language model without even looking at the data.

(c) Jacob Eisenstein 2014-2017. Work in progress.

In practice, n -gram language models tend to give perplexities in the range between 1 and V . For example, Jurafsky and Martin estimate a language model over a vocabulary of roughly 20,000 words, on 38 million tokens of text from the Wall Street Journal (Jurafsky and Martin, 2009, page 97). They report the following perplexities on a held-out set of 1.5 million tokens:

- Unigram ($n = 1$): 962
- Bigram ($n = 2$): 170
- Trigram ($n = 3$): 109

Will this trend continue?

5.5 Out-of-vocabulary words

Through this chapter, we have assumed a **closed-vocabulary** setting — the vocabulary \mathcal{V} is assumed to be a finite set. In realistic application scenarios, this assumption may not hold. Consider, for example, the problem of translating newspaper articles. The following sentence appeared in a Reuters article on January 6, 2017:⁵

The report said U.S. intelligence agencies believe Russian military intelligence, the **GRU**, used intermediaries such as **WikiLeaks**, **DCLeaks.com** and the **Guccifer 2.0** “persona” to release emails...

Suppose that you trained a language model on the Gigaword corpus,⁶ which was released in 2003. The bolded terms either did not exist at this date, or were not widely known; they are unlikely to be in the vocabulary. The same problem can occur for a variety of other terms: new technologies, previously unknown individuals, new words (e.g., *hashtag*), and numbers.

One solution is to simply mark all such terms with a special token, $\langle \text{UNK} \rangle$. While training the language model, we decide in advance on the vocabulary (often the K most common terms), and mark all other terms in the training data as $\langle \text{UNK} \rangle$. If we do not want to determine the vocabulary size in advance, an alternative approach is to simply mark the first occurrence of each word type as $\langle \text{UNK} \rangle$.

In some scenarios, we may prefer to make distinctions about the likelihood of various unknown words. This is particularly important in languages that have rich morphological

⁵Bayoumy, Y. and Strobel, W. (2017, January 6). U.S. intel report: Putin directed cyber campaign to help Trump. *Reuters*. Retrieved from <http://www.reuters.com/article/us-usa-russia-cyber-idUSKBN14Q1T8> on January 7, 2017.

⁶<https://catalog.ldc.upenn.edu/LDC2003T05>

systems, with many inflections for each word. For example, Spanish is only moderately complex from a morphological perspective, yet each verb has dozens of inflected forms. In such languages, there will necessarily be many word types that we do not encounter in a corpus, which are nonetheless predictable from the morphological rules of the language. To use a somewhat contrived English example, if *transfenestrate* is in the vocabulary, our language model should assign a non-zero probability to the past tense *transfenestrated*, even if it does not appear in the training data.

One way to accomplish this is to supplement word-level language models with **character-level language models**. Such models can use n -grams or RNNs, but with a fixed vocabulary equal to the set of ASCII or Unicode characters. For example Ling et al. (2015b) propose an LSTM model over characters, and Kim (2014) employ a **convolutional neural network** (LeCun and Bengio, 1995). A more linguistically motivated approach is to segment words into meaningful subword units, known as **morphemes** (see chapter 8). For example, Botha and Blunsom (2014) induce vector representations for morphemes, which they build into a log-bilinear language model; Bhatia et al. (2016) incorporate morpheme vectors into an LSTM.

Exercises

1. exercises tk

Chapter 6

Sequence labeling

In sequence labeling, we want to assign tags to words, or more generally, we want to assign discrete labels to elements in a sequence. There are many applications of sequence labeling in natural language processing, and chapter 7 presents an overview. One of the most classic application of sequence labeling is **part-of-speech tagging**, which involves tagging each word by its grammatical category. Coarse-grained grammatical categories include **NOUNs**, which describe things, properties, or ideas, and **VERBs**, which describe actions and events. Given a simple sentence like,

(6.1) They can fish.

we would like to produce the tag sequence N V V, with the modal verb *can* labeled as a verb in this simplified example.

6.1 Sequence labeling as classification

One way to solve tagging problems is to treat them as classification. We can write $f((\mathbf{w}, m), y)$ to indicate the feature function for applying tag y to word w_m in the sequence w_1, w_2, \dots, w_M . A simple tagging model would have a single base feature, the word itself:

$$f((\mathbf{w} = \text{they can fish}, m = 1), N) = \langle \text{they}, N \rangle \quad [6.1]$$

$$f((\mathbf{w} = \text{they can fish}, m = 2), V) = \langle \text{can}, V \rangle \quad [6.2]$$

$$f((\mathbf{w} = \text{they can fish}, m = 3), V) = \langle \text{fish}, V \rangle. \quad [6.3]$$

Here the feature function takes three arguments as input: the sentence to be tagged (*they can fish* in all cases), the proposed tag (e.g., N or V), and the word token to which this tag is applied. This simple feature function then returns a single feature: a tuple including the word to be tagged and the tag that has been proposed. If the vocabulary size is V

and the number of tags is K , then there are $V \times K$ features. Each of these features must be assigned a weight. These weights can be learned from a labeled dataset using a classification algorithm such as perceptron, but this isn't necessary in this case: it would be equivalent to define the classification weights directly, with $\theta_{w,y} = 1$ for the tag y most frequently associated with word w , and $\theta_{w,y} = 0$ for all other tags.

However, it is easy to see that this simple classification approach can go wrong. Consider the word *fish*, which often describes an animal rather than an activity; in these cases, *fish* should be tagged as a noun. To tag ambiguous words correctly, the tagger must rely on context, such as the surrounding words. We can build this context into the feature set by incorporating the surrounding words as additional features:

$$\begin{aligned} f((w = \text{they can fish}, 1), N) = \{ & \langle w_i = \text{they}, y_i = N \rangle, \\ & \langle w_{i-1} = \diamond, y_i = N \rangle, \\ & \langle w_{i+1} = \text{can}, y_i = N \rangle \} \end{aligned} \quad [6.4]$$

$$\begin{aligned} f((w = \text{they can fish}, 2), V) = \{ & \langle w_i = \text{can}, y_i = V \rangle, \\ & \langle w_{i-1} = \text{they}, y_i = V \rangle, \\ & \langle w_{i+1} = \text{fish}, y_i = V \rangle \} \end{aligned} \quad [6.5]$$

$$\begin{aligned} f((w = \text{they can fish}, 3), V) = \{ & \langle w_i = \text{fish}, y_i = V \rangle, \\ & \langle w_{i-1} = \text{can}, y_i = V \rangle, \\ & \langle w_{i+1} = \blacklozenge, y_i = V \rangle \}. \end{aligned} \quad [6.6]$$

These features contain enough information that a tagger should be able to choose the right label for the word *fish*: words that follow the modal verb *can* are likely to be verbs themselves, so the feature $\langle w_{i-1} = \text{can}, y_i = V \rangle$ should have a large positive weight.

However, even with this enhanced feature set, it may be difficult to tag some sequences correctly. One reason is that there are often relationships between the tags themselves. For example, in English it is relatively rare for a verb to follow another verb — particularly if we differentiate MODAL verbs like *can* and *should* from more typical verbs, like *give*, *transcend*, and *befuddle*. We would like to incorporate preferences **against** such tag sequences, and preferences **for** other tag sequences, such as NOUN-VERB.

The need for such preferences is best illustrated by a **garden path sentence**:

(6.2) The old man the boat.

Grammatically, the word *the* is a DETERMINER. When you read the sentence, what part of speech did you first assign to *old*? Typically, this word is an ADJECTIVE — abbreviated as J — which is a class of words that modify nouns. Similarly, *man* is usually a noun. The resulting sequence of tags is D J N D N. But this is a mistaken “garden path” interpretation, which ends up leading nowhere. It is unlikely that a determiner would directly follow a noun,¹ and particularly unlikely that the entire sentence would lack a verb. The

¹The main exception is the double object construction, as in *I gave the child a toy*.

only possible verb in the sentence is the word *man*, which can refer to the act of maintaining and piloting something — often boats. But if *man* is tagged as a verb, then *old* is seated between a determiner and a verb, and must be a noun. And indeed, adjectives can often have a second interpretation as nouns when used in this way (e.g., *the young*, *the restless*). This reasoning, in which the labeling decisions are intertwined, cannot be applied in a setting where each tag is produced by an independent classification decision.

6.2 Sequence labeling as structure prediction

As an alternative, we can think of the entire sequence of tags as a label itself. For a given sequence of words $\mathbf{w}_{1:M} = (w_1, w_2, \dots, w_M)$, there is a set of possible taggings $\mathcal{Y}(\mathbf{w}_{1:M}) = \mathcal{Y}^M$, where $\mathcal{Y} = \{N, V, D, \dots\}$ refers to the set of individual tags, and \mathcal{Y}^M refers to the set of tag sequences of length M . We can then treat the sequence labeling problem as a classification problem in the label space $\mathcal{Y}(\mathbf{w}_{1:M})$,

$$\hat{\mathbf{y}}_{1:M} = \underset{\mathbf{y}_{1:M} \in \mathcal{Y}(\mathbf{w}_{1:M})}{\operatorname{argmax}} \quad \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}_{1:M}, \mathbf{y}_{1:M}), \quad [6.7]$$

where $\mathbf{y}_{1:M} = (y_1, y_2, \dots, y_M)$ is a sequence of M tags. Note that in this formulation, we have a feature function that considers the entire tag sequence $\mathbf{y}_{1:M}$. Such a feature function can therefore include features that capture the relationships between tagging decisions, such as the preference that determiners not follow nouns, or that all sentences have verbs.

Given that the label space is exponentially large in the length of the sequence w_1, \dots, w_M , can it ever be practical to perform tagging in this way? The problem of making a series of interconnected labeling decisions is known as **inference**. Because natural language is full of interrelated grammatical structures, inference is a crucial aspect of contemporary natural language processing. In English, it is not unusual to have sentences of length $M = 20$; part-of-speech tag sets vary in size from 10 to several hundred. Taking the low end of this range, we have $|\mathcal{Y}(\mathbf{w}_{1:M})| \approx 10^{20}$, one hundred billion billion possible tag sequences. Enumerating and scoring each of these sequences would require an amount of work that is exponential in the sequence length, so inference is intractable.

However, the situation changes when we restrict the feature function. Suppose we choose features that never consider more than one tag. We can indicate this restriction as,

$$\mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, m), \quad [6.8]$$

where we use the shorthand $\mathbf{w} \triangleq \mathbf{w}_{1:M}$. The summation in Equation 6.8 means that the overall feature vector is the sum of feature vectors associated with each individual tagging decision. These features are not capable of capturing the intuitions that might help us solve garden path sentences, such as the insight that determiners rarely follow

nouns in English. But this restriction does make it possible to find the globally optimal tagging, by making a sequence of individual tagging decisions.

$$\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}) = \boldsymbol{\theta} \cdot \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, m) \quad [6.9]$$

$$= \sum_{m=1}^M \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, m) \quad [6.10]$$

$$\hat{y}_m = \operatorname{argmax}_{y_m} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, m) \quad [6.11]$$

$$\hat{\mathbf{y}} = (\hat{y}_1, \hat{y}_2, \dots, \hat{y}_M) \quad [6.12]$$

Note that we are still searching over an exponentially large set of tag sequences! But the feature set restriction results in decoupling the labeling decisions that were previous interconnected. As a result, it is not necessary to score every one of the $|\mathcal{Y}|^M$ tag sequences individually — we can find the optimal sequence by scoring the local parts of these decisions.

Now let's consider a slightly less restrictive feature function: rather than considering only individual tags, we will consider adjacent tags too. This means that we can have negative weights for infelicitous tag pairs, such as noun-determiner, and positive weights for typical tag pairs, such as determiner-noun and noun-verb. We define this feature function as,

$$\mathbf{f}(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m). \quad [6.13]$$

Let's apply this feature function to the shorter example, *they can fish*, using features for word-tag and tag-tag pairs:

$$\mathbf{f}(\mathbf{w} = \text{they can fish}, \mathbf{y} = \text{N V V}) = \sum_{m=1}^M \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) \quad [6.14]$$

$$\begin{aligned} &= \mathbf{f}(\mathbf{w}, \text{N}, \diamond, 1) \\ &\quad + \mathbf{f}(\mathbf{w}, \text{V}, \text{N}, 2) \\ &\quad + \mathbf{f}(\mathbf{w}, \text{V}, \text{V}, 3) \end{aligned} \quad [6.15]$$

$$\begin{aligned} &= \langle w_m = \text{they}, y_m = \text{N} \rangle + \langle y_m = \text{N}, y_{m-1} = \diamond \rangle \\ &\quad + \langle w_m = \text{can}, y_m = \text{V} \rangle + \langle y_m = \text{V}, y_{m-1} = \text{N} \rangle \\ &\quad + \langle w_m = \text{fish}, y_m = \text{V} \rangle + \langle y_m = \text{V}, y_{m-1} = \text{V} \rangle \\ &\quad + \langle y_m = \diamond, y_{m-1} = \text{V} \rangle. \end{aligned} \quad [6.16]$$

We end up with seven active features: one for each word-tag pair, and one for each tag-tag pair (this includes a final tag $y_{M+1} = \diamond$). These features capture what are arguably the

(c) Jacob Eisenstein 2014-2017. Work in progress.

two main sources of information for part-of-speech tagging: which tags are appropriate for each word, and which tags tend to follow each other in sequence. Given appropriate weights for these features, we can expect to make the right tagging decisions, even for difficult cases like *the old man the boat*.

The example shows that even with the restriction to the feature set shown in Equation 6.13, it is still possible to construct expressive features that are capable of solving many sequence labeling problems. But the key question is: does this restriction make it possible to perform efficient inference? The answer is yes, and the solution is the **Viterbi algorithm** (Viterbi, 1967).

6.3 The Viterbi algorithm

We now consider the inference problem,

$$\hat{y} = \underset{y}{\operatorname{argmax}} \theta \cdot f(w, y) \quad [6.17]$$

$$f(w, y) = \sum_{m=1}^M f(w, y_m, y_{m-1}, m). \quad [6.18]$$

Given this restriction on the feature function, we can solve this inference problem using **dynamic programming**, a algorithmic technique for reusing work in recurrent computations. As is often the case in dynamic programming, we begin by solving an auxiliary problem: rather than finding the best tag sequence, we simply try to compute the **score** of the best tag sequence,

$$\max_y \theta \cdot f(w, y) = \max_{y_{1:M}} \sum_{m=1}^M \theta \cdot f(w, y_m, y_{m-1}, m) \quad [6.19]$$

$$= \max_{y_{1:M}} \theta \cdot f(w, y_M, y_{M-1}, M) + \sum_{m=1}^{M-1} \theta \cdot f(w, y_m, y_{m-1}, m) \quad [6.20]$$

$$= \max_{y_M} \max_{y_{M-1}} \theta \cdot f(w, y_M, y_{M-1}, M) + \max_{y_{1:M-2}} \sum_{m=1}^{M-1} \theta \cdot f(w, y_m, y_{m-1}, m). \quad [6.21]$$

In this derivation, we first removed the final element $\theta \cdot f(w, y_M, y_{M-1}, M)$ from the sum over the sequence, and then we adjusted the scope of the the max operation, since the elements $(y_1 \dots y_{M-2})$ are irrelevant to the final term.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Let us now define the **Viterbi variable**,

$$v_m(k) \triangleq \max_{\mathbf{y}_{1:m-1}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + \sum_{n=1}^{m-1} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n), \quad [6.22]$$

where lower-case m indicates any position in the sequence, and $k \in \mathcal{Y}$ indicates a tag for that position. The variable $v_m(k)$ represents the score of the best tag sequence $(\hat{y}_1, \hat{y}_2, \dots, \hat{y}_m)$ that terminates in $\hat{y}_m = k$. From this definition, we can compute the score of the best tagging of the sequence by plugging the Viterbi variables $v_M(\cdot)$ into Equation 6.21,

$$\max_{\mathbf{y}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}) = \max_k v_M(k). \quad [6.23]$$

Now, let us look more closely at how we can compute these Viterbi variables.

$$v_m(k) \triangleq \max_{\mathbf{y}_{1:m-1}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + \sum_{n=1}^{m-1} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n) \quad [6.24]$$

$$\begin{aligned} &= \max_{y_{m-1}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) \\ &\quad + \max_{\mathbf{y}_{1:m-2}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_{m-1}, y_{m-2}) + \sum_{n=1}^{m-2} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n) \end{aligned} \quad [6.25]$$

$$= \max_{y_{m-1}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + v_{m-1}(y_{m-1}) \quad [6.26]$$

$$v_1(y) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y, \diamond, 1). \quad [6.27]$$

Equation 6.26 is a **recurrence** for computing the Viterbi variables: each $v_m(k)$ can be computed in terms of $v_{m-1}(\cdot)$, and so on. We can therefore step forward through the sequence, computing first all variables $v_1(\cdot)$ from Equation 6.27, and then computing all variables $v_2(\cdot)$, $v_3(\cdot)$, and so on, until we reach the final set of variables $v_M(\cdot)$.

Graphically, it is customary to arrange these variables in a matrix, with the sequence index m on the columns, and the tag index k on the rows. In this representation, each $v_{m-1}(k)$ is connected to each $v_m(k')$, forming a **trellis**, as shown in Figure 6.1. As shown in the figure, special nodes are set aside for the start and end states.

Our real goal is to find the best scoring sequence, not simply to compute its score. But as is often the case in dynamic programming, solving the auxiliary problem gets us almost all the way to our original goal. Recall that each $v_m(k)$ represents the score of the best tag sequence ending in that tag k in position m . To compute this, we maximize over possible values of y_{m-1} . If we keep track of the tag that maximizes this choice at each step, then we can walk backwards from the final tag, and recover the optimal tag sequence. This is indicated in Figure 6.1 by the solid blue lines, which we trace back from the final position.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 10 The Viterbi algorithm.

```

for  $k \in \{0, \dots, K\}$  do
   $v_1(k) = \theta \cdot \mathbf{f}(\mathbf{w}, k, \diamond, m)$ 
for  $m \in \{2, \dots, M\}$  do
  for  $k \in \{0, \dots, K\}$  do
     $v_m(k) = \max_{k'} \theta \cdot \mathbf{f}(\mathbf{w}, k, k', m) + v_{m-1}(k')$ 
     $b_m(k) = \operatorname{argmax}_{k'} \theta \cdot \mathbf{f}(\mathbf{w}, k, k', m) + v_{m-1}(k')$ 
 $y_M = \operatorname{argmax}_k v_M(k) + \theta \cdot \mathbf{f}(\mathbf{w}, \diamond, k, M + 1)$ 
for  $m \in \{M - 1, \dots, 1\}$  do
   $y_m = b_m(y_{m+1})$ 
return  $y_{1:M}$ 
  
```

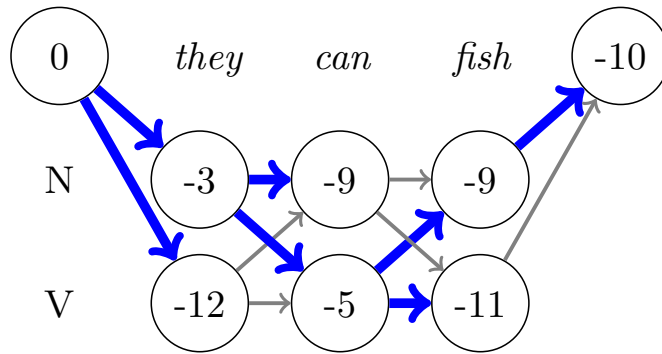


Figure 6.1: The trellis representation of the Viterbi variables, for the example *they can fish*, using the weights shown in Table 6.1.

These “back-pointers” are written $b_m(k)$, indicating the optimal tag y_{m-1} on the path to $Y_m = k$.

Why does this work? We can make an inductive argument. Suppose k is indeed the optimal tag for word m , and we now need to decide on the tag y_{m-1} . Because we make the inductive assumption that we know $y_m = k$, and because the feature function is restricted to adjacent tags, we need not consider any of the tags $y_{m+1:M}$; these tags, and the features that describe them, are irrelevant to the inference of y_{m-1} , given that we have $y_m = k$. Thus, we are looking for the tag \hat{y}_{m-1} that maximizes,

$$\hat{y}_{m-1} = \operatorname{argmax}_{y_{m-1}} \theta \cdot \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + \max_{\mathbf{y}_{1:m-2}} \sum_{n=1}^{m-1} \theta \cdot \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n) \quad [6.28]$$

$$= \operatorname{argmax}_{y_{m-1}} \theta \cdot \mathbf{f}(\mathbf{w}, k, y_{m-1}, m) + v_{m-1}(y_{m-1}), \quad [6.29]$$

	<i>they</i>	<i>can</i>	<i>fish</i>
N	-2	-3	-3
V	-10	-1	-3

(a) Weights for emission features.

	N	V	◆
◇	-1	-2	$-\infty$
N	-3	-1	-12
V	-1	-3	-1

(b) Weights for transition features. The “from” tags are on the columns, and the “to” tags are on the rows.

Table 6.1: Feature weights for the example trellis shown in Figure 6.1. Emission weights from \diamond and \blacklozenge are implicitly set to $-\infty$.

which we obtain by plugging in the definition of the Viterbi variable. The value \hat{y}_{m-1} was identified during forward pass, when computing the value of the Viterbi variable $v_m(k)$.

The complete Viterbi algorithm is shown in Algorithm 10. This formalizes the recurrences that were described in the previous paragraphs, and handles the boundary conditions at the start and end of the sequence. Specifically, when computing the initial Viterbi variables $v_1(\cdot)$, we use a special tag, \diamond , to indicate the start of the sequence. When computing the final tag Y_M , we use another special tag, \blacklozenge , to indicate the end of the sequence. These special tags enable the use of transition features for the tags that begin and end the sequence: for example, conjunctions are unlikely to end sentences in English, so we would like a large negative weight for the feature $\langle \text{CC}, \blacklozenge \rangle$; nouns are relatively likely to appear at the beginning of sentences, so we would like a more positive (or less negative) weight for the feature $\langle \diamond, \text{N} \rangle$.

What is the complexity of this algorithm? If there are K tags and M positions in the sequence, then there are $M \times K$ Viterbi variables to compute. Computing each variable requires finding a maximum over K possible predecessor tags. The total computation cost of populating the trellis is therefore $\mathcal{O}(MK^2)$, with an additional factor for the number of active features at each position. After completing the trellis, we simply trace the backwards pointers to the beginning of the sequence, which takes $\mathcal{O}(M)$ operations.

6.3.1 Example

To illustrate the Viterbi algorithm with an example, let us consider the minimal tagset $\{\text{N}, \text{V}\}$, corresponding to nouns and verbs. Even in this tagset, there is considerable ambiguity: for example, the words *can* and *fish* can each take both tags. Of the $2 \times 2 \times 2 = 8$ possible taggings for the sentence *they can fish*, four are possible given these possible tags, and two are grammatical. (The tagging *they*/N *can*/V *fish*/N corresponds to the scenario of putting fish into cans.)

To begin, we use the feature weights defined in Table 6.1. These weights are used to

(c) Jacob Eisenstein 2014-2017. Work in progress.

incrementally fill in the trellis. As described in Algorithm 10, we fill in the cells from left to right, with each column corresponding to a word in the sequence. As we fill in the cells, we must keep track of the back-pointers $b_m(k)$ — the previous cell that maximizes the score of tag k at word m . These are represented in the figure with the thick blue lines. At the end of the algorithm, we recover the optimal tag sequence by tracing back the optimal path from the final position, $(M + 1, \blacklozenge)$.

6.3.2 Higher-order features

The Viterbi algorithm was made possible by a restriction of the features to consider only pairs of adjacent tags. In a sense, we can think of this as a bigram language model, at the tag level. A natural question is how to generalize Viterbi to tag trigrams, which would involve the following feature decomposition:

$$f(\mathbf{w}, \mathbf{y}) = \sum_m^M f(\mathbf{w}, y_m, y_{m-1}, y_{m-2}, m). \quad [6.30]$$

One possibility is to take the Cartesian product of the tagset with itself, $\mathcal{Y}^{(2)} = \mathcal{Y} \times \mathcal{Y}$. The tags in this product space are ordered pairs, representing adjacent tags at the token level: for example, the tag $\langle N, V \rangle$ would represent a noun followed by a verb. Transitions between such tags must be consistent: we can have a transition from $\langle N, V \rangle$ to $\langle V, N \rangle$ (corresponding to the token-level tag sequence $N V N$), but not from $\langle N, V \rangle$ to $\langle N, N \rangle$, which would not correspond to any token-level tag sequence. This constraint can be enforced in the feature weights, with $\theta_{\langle\langle a,b \rangle, \langle c,d \rangle\rangle} = -\infty$ if $b \neq c$. The remaining feature weights can encode preferences for and against various tag trigrams.

In the Cartesian product tag space, there are K^2 tags, suggesting that the time complexity will increase to $\mathcal{O}(MK^4)$. However, it is unnecessary to max over predecessor tag bigrams that are incompatible with the current tag bigram. By exploiting these constraints, it is possible to limit the time complexity to $\mathcal{O}(MK^3)$. The space complexity is $\mathcal{O}(MK^2)$. In general, the time and space complexity of higher-order Viterbi grows exponentially with the order of the tag n -grams that are considered in the feature decomposition.

6.4 Hidden Markov Models

We now consider how to learn the weights θ that parametrize the Viterbi sequence labeling algorithm. We begin with a probabilistic approach. Recall that the probabilistic Naïve Bayes classifier selects the label y to maximize $p(y | \mathbf{x}) \propto p(y, \mathbf{x})$. In probabilistic sequence labeling, our goal is similar: select the tag sequence that maximizes $p(\mathbf{y} | \mathbf{w}) \propto p(\mathbf{y}, \mathbf{w})$. Just as Naïve Bayes could be cast as a linear classifier maximizing $\theta \cdot \mathbf{f}(\mathbf{x}, y)$, we can cast

our probabilistic classifier as a linear decision rule. Furthermore, the feature restriction in Equation 6.13 can be viewed as a conditional independence assumption on the random variables \mathbf{y} . Thanks to this assumption, it is possible to perform inference using the Viterbi algorithm.

Naïve Bayes was introduced as a generative model — a probabilistic story that explains the observed data as well as the hidden label. A similar story can be constructed for probabilistic sequence labeling: first, we draw the tags from a prior distribution, $\mathbf{y} \sim p(\mathbf{y})$; next, we draw the tokens from a conditional likelihood distribution, $\mathbf{w} \mid \mathbf{y} \sim p(\mathbf{w} \mid \mathbf{y})$. However, for inference to be tractable, additional independence assumptions are required. Here we make two assumptions. First, the probability of each token depends only on its tag, and not on any other element in the sequence:

$$p(\mathbf{w} \mid \mathbf{y}) = \prod_{m=1}^M p(w_m \mid y_m). \quad [6.31]$$

Next, we introduce an independence assumption on the form of the prior distribution over labels: each label y_m depends only on its predecessor,

$$p(\mathbf{y}) = \prod_{m=1}^M p(y_m \mid y_{m-1}), \quad [6.32]$$

where $y_0 = \diamond$ in all cases. Due to this **Markov** assumption, probabilistic sequence labeling models are known as **hidden Markov models** (HMMs). We now state the generative model under these independence assumptions,

- For $m \in (1, 2, \dots, M)$,
 - draw $y_m \mid y_{m-1} \sim \text{Categorical}(\boldsymbol{\lambda}_{y_{m-1}})$;
 - draw $w_m \mid y_m \sim \text{Categorical}(\phi_{y_m})$

This generative story formalizes the hidden Markov model. Given the parameters $\boldsymbol{\lambda}$ and ϕ , we can compute $p(\mathbf{w}, \mathbf{y})$ for any token sequence \mathbf{w} and tag sequence \mathbf{y} . The HMM is often represented as a **graphical model** (Wainwright and Jordan, 2008), as shown in Figure 6.2. This representation makes the independence assumptions explicit: if a variable v_1 is probabilistically conditioned on another variable v_2 , then there is an arrow $v_2 \rightarrow v_1$ in the diagram. If there are no arrows between v_1 and v_2 , they are **conditionally independent**, given each variable’s **Markov blanket**. In the hidden Markov model, the Markov blanket for each tag y_m includes the “parent” y_{m-1} , and the “children” y_{m+1} and w_m .²

²In general graphical models, a variable’s Markov blanket includes its parents, children, and its children’s other parents (Murphy, 2012).

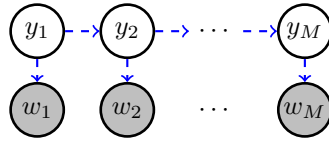


Figure 6.2: Graphical representation of the hidden Markov model. Arrows indicate probabilistic dependencies.

It is important to reflect on the implications of the HMM independence assumptions. A non-adjacent pair of tags y_m and y_n are conditionally independent; if $m < n$ and we are given y_{n-1} , then y_m offers no additional information about y_n . However, if we are not given any information about the tags in a sequence, then all tags are probabilistically coupled.

6.4.1 Estimation

The hidden Markov model has two groups of parameters:

Emission probabilities. The probability $p_e(w_m | y_m; \phi)$ is the emission probability, since the words are treated as probabilistically “emitted”, conditioned on the tags.

Transition probabilities. 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.

Both of these groups of parameters are typically computed from relative frequency estimation on a labeled corpus,

$$\phi_{k,i} \triangleq \Pr(W_m = i | Y_m = k) = \frac{\text{count}(W_m = i, Y_m = k)}{\text{count}(Y_m = k)}$$

$$\lambda_{k,k'} \triangleq \Pr(Y_m = k' | Y_{m-1} = k) = \frac{\text{count}(Y_m = k', Y_{m-1} = k)}{\text{count}(Y_{m-1} = k)}.$$

Smoothing is more important for the emission probability than the transition probability, because the event space is much larger. Smoothing techniques such as additive smoothing, interpolation, and backoff (see chapter 5) can all be applied here.

6.4.2 Inference

The goal of inference in the hidden Markov model is to find the highest probability tag sequence,

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} p(\mathbf{y} | \mathbf{w}). \quad [6.33]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

As in Naïve Bayes, it is equivalent to find the tag sequence with the highest **log**-probability, since the log function is monotonically increasing. It is furthermore equivalent to maximize the joint probability $p(\mathbf{y}, \mathbf{w}) = p(\mathbf{y} | \mathbf{w}) \times p(\mathbf{w}) \propto p(\mathbf{y} | \mathbf{w})$, which is proportional to the conditional probability. Therefore, we can reformulate the inference problem as,

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \log p(\mathbf{y}, \mathbf{w}). \quad [6.34]$$

We can now apply the HMM independence assumptions:

$$\log p(\mathbf{y}, \mathbf{w}) = \log p(\mathbf{y}) + \log p(\mathbf{w} | \mathbf{y}) \quad [6.35]$$

$$= \sum_{m=1}^M \log p_y(y_m | y_{m-1}) + \log p_{w|y}(w_m | y_m) \quad [6.36]$$

$$= \sum_{m=1}^M \log \lambda_{y_m, y_{m-1}} + \log \phi_{y_m, w_m}. \quad [6.37]$$

This log probability can be rewritten as a dot product of weights and features,

$$\log p(\mathbf{y}, \mathbf{w}) = \sum_{m=1}^M \log \lambda_{y_m, y_{m-1}} + \log \phi_{y_m, w_m} \quad [6.38]$$

$$= \sum_{m=1}^M \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m), \quad [6.39]$$

where the feature function is defined,

$$\mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) = \{\langle y_m, y_{m-1} \rangle, \langle y_m, w_m \rangle\}, \quad [6.40]$$

and the weight vector $\boldsymbol{\theta}$ encodes the log-parameters $\log \lambda$ and $\log \phi$.

This derivation shows that HMM inference can be viewed as an application of the Viterbi decoding algorithm, given an appropriately defined feature function and weight vector. The local product $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$ can be interpreted probabilistically,

$$\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) = \log p_y(y_m | y_{m-1}) + \log p_{w|y}(w_m | y_m) \quad [6.41]$$

$$= \log p(y_m, w_m | y_{m-1}). \quad [6.42]$$

Now recall the definition of the Viterbi variables,

$$v_m(k) = \max_{\mathbf{y}_{1:m-1}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, Y_m = k, y_{m-1}, m) + \sum_{n=1}^{m-1} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_n, y_{n-1}, n) \quad [6.43]$$

$$= \max_{\mathbf{y}_{1:m-1}} \log p_{y_m, w_m | y_{m-1}}(k, w_m | y_{m-1}) + \sum_{n=1}^{m-1} \log p(y_n, w_n | y_{n-1}) \quad [6.44]$$

$$= \max_{\mathbf{y}_{1:m-1}} \log p(\mathbf{y}_{1:m-1}, Y_m = k, \mathbf{w}_{1:m}). \quad [6.45]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

In words, the Viterbi variable $v_m(k)$ is the log probability of the best tag sequence ending in $Y_m = k$, joint with the word sequence $\mathbf{w}_{1:m}$. The log probability of the best complete tag sequence is therefore,

$$\max_{\mathbf{y}_{1:M}} \log p(\mathbf{y}_{1:M}, \mathbf{w}_{1:M}) = \max_{y_M} \log p_y(\diamond | y_M) + v_M(y_M). \quad [6.46]$$

The Viterbi algorithm can also be implemented using probabilities, rather than log probabilities. In this case, each $v_m(k)$ is equal to,

$$v_m(k) = \max_{\mathbf{y}_{1:m-1}} p(\mathbf{y}_{1:m-1}, Y_m = k, \mathbf{w}_{1:m}) \quad [6.47]$$

$$= \max_{y_{m-1}} p(Y_m = k, w_m | y_{m-1}) \times \max_{\mathbf{y}_{1:m-2}} p(\mathbf{y}_{1:m-2}, y_{m-1}, \mathbf{w}_{1:m-1}) \quad [6.48]$$

$$= \max_{y_{m-1}} p(Y_m = k, w_m | y_{m-1}) \times v_{m-1}(y_{m-1}) \quad [6.49]$$

$$= p_E(w_m | Y_m = k) \times \max_{y_{m-1}} p_T(y_m | y_{m-1}) \times v_{m-1}(y_{m-1}) \quad [6.50]$$

$$= \max_{y_{m-1}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, Y_m = k, y_{m-1}, m)) \times v_{m-1}(y_{m-1}). \quad [6.51]$$

In the final line, we use the fact that $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m) = \log p(y_m, w_m | y_{m-1})$, and exponentiate the dot product to obtain the probability.

In practice, the probabilities tend towards zero over long sequences, so the log-probability version of Viterbi is more practical from the standpoint of numerical stability. However, this version connects to a broader literature on inference in graphical models. Each Viterbi variable is computed by **maximizing** over a set of **products**. Thus, the Viterbi algorithm is a special case of the **max-product algorithm** for inference in graphical models (Wainwright and Jordan, 2008).

6.4.3 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 $\{\mathbf{w}^{(1)}, \mathbf{w}^{(2)}, \dots, \mathbf{w}^{(N)}\}$, estimate the transition and emission distributions.

The Viterbi algorithm solves the decoding problem. We'll talk about the learning problem in § 6.6. Let's now consider how to compute the marginal likelihood $p(\mathbf{w}) = \sum_{\mathbf{y}} p(\mathbf{w}, \mathbf{y})$, which involves summing over all possible tag sequences. There are at least two reasons we might want to do this:

(c) Jacob Eisenstein 2014-2017. Work in progress.

Language modeling Note that the probability $p(\mathbf{w})$ is also computed by the language models that were discussed in chapter 5. In those language models, we used only unlabeled corpora, conditioning each token w_m on previous tokens. An HMM-based language model would leverage a corpus of part-of-speech annotations, and therefore might be expected to generalize better than an n-gram language model — for example, it would be more likely to assign positive probability to a nonsense grammatical sentence like *colorless green ideas sleep furiously*.

Tag marginals It is often important to compute marginal probabilities of individual tags, $p(y_m \mid \mathbf{w}_{1:M})$. This is the probability distribution over tags for token m , conditioned on all of the words $\mathbf{w}_{1:M}$. For example, we might like to know the probability that a given word is tagged as a verb, regardless of how all the other words are tagged. We will discuss how to compute this probability in § 6.5.3, but as a preview, we will use the following form,

$$p(y_m \mid \mathbf{w}_{1:M}) = \frac{p(y_m, \mathbf{w}_{1:M})}{p(\mathbf{w}_{1:M})}, \quad [6.52]$$

which involves the marginal likelihood in the denominator.

We can compute the marginal likelihood using a dynamic program that is nearly identical to the Viterbi algorithm. We will use probabilities for now, and show the conversion to log-probabilities later. The core of the algorithm is to compute a set of **forward variables**,

$$\alpha_m(k) \triangleq p(Y_m = k, \mathbf{w}_{1:m}). \quad [6.53]$$

From this definition, we can compute the marginal likelihood by summing over the final forward variables,

$$p(\mathbf{w}) = \sum_{k \in \mathcal{Y}} p(Y_M = k, \mathbf{w}_{1:M}) \quad [6.54]$$

$$= \sum_{k \in \mathcal{Y}} \alpha_M(k). \quad [6.55]$$

To capture the probability of terminating the sequence on each possible tag Y_M , we can pad the end of \mathbf{w} with an extra token ■, which can only be emitted from the stop tag ♦.

(c) Jacob Eisenstein 2014-2017. Work in progress.

The forward variables themselves can be computed recursively,

$$\alpha_m(k) = \mathbf{p}(Y_m = k, \mathbf{w}_{1:m}) \quad [6.56]$$

$$= \mathbf{p}(w_m | Y_m = k) \times \Pr(Y_m = k, \mathbf{w}_{1:m-1}) \quad [6.57]$$

$$= \mathbf{p}(w_m | Y_m = k) \times \sum_{k' \in \mathcal{Y}} \Pr(Y_m = k, Y_{m-1} = k', \mathbf{w}_{1:m-1}) \quad [6.58]$$

$$= \mathbf{p}(w_m | Y_m = k) \times \sum_{k' \in \mathcal{Y}} \Pr(Y_m = k | Y_{m-1} = k') \times \Pr(Y_{m-1} = k', \mathbf{w}_{1:m-1}) \quad [6.59]$$

$$= \mathbf{p}(w_m | Y_m = k) \times \sum_{k' \in \mathcal{Y}} \Pr(Y_m = k | Y_{m-1} = k') \times \alpha_{m-1}(k') \quad [6.60]$$

$$= \sum_{k' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}_{1:M}, Y_m = k, Y_{m-1} = k', m)) \times \alpha_{m-1}(k'). \quad [6.61]$$

The derivation relies on the independence assumptions in the hidden Markov model: W_m depends only on Y_m , and Y_m is conditionally independent from $W_{1:m-1}$ and all tags, given Y_{m-1} . We complete the derivation by introducing Y_{m-1} and summing over all possible values, and by then applying the chain rule to obtain the final recursive form.

Procedurally, we compute the forward variables in just the same way as we compute the Viterbi variables: we first compute all $\alpha_1(\cdot)$, then all $\alpha_2(\cdot)$, and so on. We initialize each $\alpha_0(k) = \mathbf{p}(Y_m = k | Y_{m-1} = \diamond)$, to capture the transition probability from the start symbol. Comparing Equation 6.60 to Equation 6.50, the sole difference is that instead of maximizing over possible values of Y_{m-1} , we sum. Just as the Viterbi algorithm is a special case of the max-product algorithm for inference in graphical models, the forward algorithm is a special case of the **sum-product** algorithm for computing marginal likelihoods.

In practice, it is numerically more stable to compute the marginal log-probability. In the log domain, the forward recurrence is,

$$\alpha_m(k) \triangleq \log \mathbf{p}(\mathbf{w}_{1:m}, Y_m = k) \quad [6.62]$$

$$\begin{aligned} &= \log \sum_{k' \in \mathcal{Y}} \exp(\log \mathbf{p}(w_m | Y_m = k) + \log \Pr(Y_m = k | Y_{m-1} = k')) \\ &\quad + \log \mathbf{p}(\mathbf{w}_{1:m-1}, Y_{m-1} = k')) \end{aligned} \quad [6.63]$$

$$= \log \sum_{k' \in \mathcal{Y}} \exp(\log \mathbf{p}(w_m | Y_m = k) + \log \Pr(Y_m = k | Y_{m-1} = k') + \alpha_{m-1}(k')) \quad [6.64]$$

$$= \log \sum_{k' \in \mathcal{Y}} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}_{1:M}, Y_m = k, Y_{m-1} = k', m) + \alpha_{m-1}(k')). \quad [6.65]$$

Scientific programming libraries provide numerically robust implementations of the log-sum-exp function, which should prevent overflow and underflow from exponentiation.

(c) Jacob Eisenstein 2014-2017. Work in progress.

6.4.4 *Semiring Notation and the Generalized Viterbi Algorithm

We have now seen the Viterbi and Forward recurrences, each of which can be performed over probabilities or log probabilities. These four recurrences are closely related, and can in fact be expressed as a single recurrence in a more general notation, known as **semiring algebra**. We use the symbol \oplus to represent generalized addition, and the symbol \otimes to represent generalized multiplication.³ Given these operators, we can denote a generalized Viterbi recurrence as,

$$v_m(k) = \bigoplus_{k' \in \mathcal{Y}} \theta \cdot f(w, Y_m = k, Y_{m-1} = k', m) \otimes v_{m-1}(k'). \quad [6.66]$$

Each recurrence that we have seen so far is a special case of this generalized Viterbi recurrence:

- In the max-product Viterbi recurrence over probabilities, the \oplus operation corresponds to maximization, and the \otimes operation corresponds to multiplication.
- In the forward recurrence over probabilities, the \oplus operation corresponds to addition, and the \otimes operation corresponds to multiplication.
- In the max-product Viterbi recurrence over log-probabilities, the \oplus operation corresponds to maximization, and the \otimes operation corresponds to addition. (This is sometimes called the **tropical semiring**, in honor of the Brazilian mathematician Imre Simon.)
- In the forward recurrence over log-probabilities, the \oplus operation corresponds to log-addition, $a \oplus b = \log(e^a + e^b)$. The \otimes operation corresponds to addition.

The mathematical abstraction offered by semiring notation can be applied to the software implementations of these algorithms, yielding concise and modular implementations. The OPENFST library (Allauzen et al., 2007) is an example of a software package in which the algorithms are parametrized by the choice of semiring.

6.5 Discriminative sequence labeling

Today, hidden Markov models are rarely used for supervised sequence labeling. This is because HMMs are limited to only two phenomena:

- Word-tag probabilities, via the emission probability $p_E(w_m | y_n)$;

³In a semiring, the addition and multiplication operators must both obey associativity, and multiplication must distribute across addition; the addition operator must be commutative; there must be additive and multiplicative identities $\bar{0}$ and $\bar{1}$, such that $a \oplus \bar{0} = a$ and $a \otimes \bar{1} = a$; and there must be a multiplicative annihilator $\bar{0}$, such that $a \otimes \bar{0} = \bar{0}$.

- local context, via the transition probability $p_T(y_m \mid y_{m-1})$.

However, as we have seen, the Viterbi algorithm can be applied to much more general feature sets, as long as the decomposition $f(\mathbf{w}, \mathbf{y}) = \sum_{m=1}^M f(\mathbf{w}, y_m, y_{m-1}, m)$ is observed. In this section, we discuss methods for learning the weights on such features. However, let's first pause to ask what additional features might be needed.

Word affix features. Consider the problem of part-of-speech tagging on the first four lines of the poem *Jabberwocky* (Carroll, 1917):

(6.3) 'Twas brillig, and the slithy toves
Did gyre and gimble in the wabe:
All mimsy were the borogoves,
And the mome raths outgrabe.

Many of these words are made up, so you would have no information about their probabilities of being associated with any particular part of speech. Yet it is not so hard to see what their grammatical roles might be in this passage. Context helps: for example, the word *slithy* follows the determiner *the*, and therefore is likely to be a noun or adjective. Which do you think is more likely? The suffix *-thy* is found in a number of adjectives — e.g., *frothy, healthy, pithy, worthy*. The suffix is also found in a handful of nouns — e.g., *apathy, sympathy* — but nearly all of these nouns contain *-pathy*, unlike *slithy*. The suffix gives some evidence that *slithy* is an adjective, and indeed it is: later in the text we find that it is a combination of the adjectives *lithe* and *slimy*.⁴

Fine-grained context. Another useful source of information is fine-grained context — that is, contextual information that is more specific than the previous tag. For example, consider the noun phrases *this fish* and *these fish*. Many part-of-speech tagsets distinguish between singular and plural nouns, but do not distinguish between singular and plural determiners; for example, the Penn Treebank tagset follows these conventions. A hidden Markov model would be unable to correctly label *fish* as singular or plural in both of these cases, because it only has access to two features: the preceding tag (determiner in both cases) and the word (*fish* in both cases). The classification-based tagger discussed in § 6.1 had the ability to use preceding and succeeding words as features, and we would like to incorporate this information into a sequence labeling algorithm.

⁴**Morphology** is the study of how words are formed from smaller linguistic units. Computational approaches to morphological analysis are touched on in chapter 8; Bender (2013) provides a good overview of the underlying linguistic principles.

Example Suppose we have the tagging D J N (determiner, adjective, noun) for the sequence *the slithy toves* in Jabberwocky, so that

$$\begin{aligned} \mathbf{w} &= \text{the slithy toves} \\ \mathbf{y} &= \text{D J N}. \end{aligned}$$

We now create the feature vector for this example, assuming that we have word-tag features (indicated by prefix W), tag-tag features (indicated by prefix T), and suffix features (indicated by prefix M). We assume access to a method for extracting the suffix *-thy* from *slithy*, *-es* from *toves*, and \emptyset from *the*, indicating that this word has no suffix. The resulting feature vector is,

$$\begin{aligned} \mathbf{f}(\text{the slithy toves}, \text{D J N}) = & \{ \langle W : \text{the}, \text{D} \rangle, \langle M : \emptyset, \text{D} \rangle, \langle T : \diamond, \text{D} \rangle \\ & \langle W : \text{slithy}, \text{J} \rangle, \langle M : \text{-thy}, \text{J} \rangle, \langle T : \text{D}, \text{J} \rangle \\ & \langle W : \text{toves}, \text{N} \rangle, \langle M : \text{-es}, \text{N} \rangle, \langle T : \text{J}, \text{N} \rangle \\ & \langle T : \text{N}, \diamond \rangle \}. \end{aligned}$$

We now consider several discriminative methods for learning feature weights in sequence labeling. In § 1.2, we considered three types of discriminative classifiers: perceptron, support vector machine, and logistic regression. Each of these classifiers has a structured equivalent, enabling it to be trained from labeled sequences rather than individual tokens.

6.5.1 Structured perceptron

The perceptron classifier updates its weights by increasing the weights for features that are associated with the correct label, and decreasing the weights for features that are associated with incorrectly predicted labels:

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{x}, y) \quad [6.67]$$

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} + \mathbf{f}(\mathbf{x}, y) - \mathbf{f}(\mathbf{x}, \hat{y}). \quad [6.68]$$

We can apply exactly the same update in the case of structure prediction,

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}) \quad [6.69]$$

$$\boldsymbol{\theta}^{(t+1)} \leftarrow \boldsymbol{\theta}^{(t)} + \mathbf{f}(\mathbf{w}, \mathbf{y}) - \mathbf{f}(\mathbf{w}, \hat{\mathbf{y}}). \quad [6.70]$$

This learning algorithm is called **structured perceptron**, because it learns to predict the structured output \mathbf{y} . The key difference is that instead of computing \hat{y} by enumerating

the entire set \mathcal{Y} , we use the Viterbi algorithm to search this set efficiently. In this case, the output structure is the sequence of tags (y_1, y_2, \dots, y_M) ; the algorithm can be applied to other structured outputs as long as efficient inference is possible. As in perceptron classification, weight averaging is crucial to get good performance (see § 1.2.2).

Example For the example *They can fish*, suppose the reference tag sequence is N V V, but our tagger incorrectly returns the tag sequence N V N. Given **feature templates** $\langle w_m, y_m \rangle$ and $\langle y_{m-1}, y_m \rangle$, the corresponding structured perceptron update is:

$$\theta_{\langle \text{fish}, V \rangle} \leftarrow \theta_{\langle \text{fish}, V \rangle} + 1 \quad [6.71]$$

$$\theta_{\langle \text{fish}, N \rangle} \leftarrow \theta_{\langle \text{fish}, N \rangle} - 1 \quad [6.72]$$

$$\theta_{\langle V, V \rangle} \leftarrow \theta_{\langle V, V \rangle} + 1 \quad [6.73]$$

$$\theta_{\langle V, N \rangle} \leftarrow \theta_{\langle V, N \rangle} - 1 \quad [6.74]$$

$$\theta_{\langle V, \blacklozenge \rangle} \leftarrow \theta_{\langle V, \blacklozenge \rangle} + 1 \quad [6.75]$$

$$\theta_{\langle N, \blacklozenge \rangle} \leftarrow \theta_{\langle N, \blacklozenge \rangle} - 1. \quad [6.76]$$

6.5.2 Structured Support Vector Machines

Large-margin classifiers such as the support vector machine improve on the perceptron by learning weights that push the classification boundary away from the training instances. In many cases, large-margin classifiers outperform the perceptron, so we would like to apply similar ideas to sequence labeling. A support vector machine in which the output is a structured object, such as a sequence, is called a **structured support vector machine** (Tsochantaridis et al., 2004).⁵

In classification, we formalized the large-margin constraint as,

$$\forall y \neq y^{(i)}, \theta \cdot \mathbf{f}(x, y^{(i)}) - \theta \cdot \mathbf{f}(x, y) \geq 1, \quad [6.77]$$

which says that we require a margin of at least 1 between the scores for all labels y that are not equal to the correct label $y^{(i)}$. The weights θ are then learned by constrained optimization (see § 1.3.2).

We can apply this idea to sequence labeling by formulating an equivalent set of constraints for all possible labelings $\mathcal{Y}(w)$ for an input w . However, there are two problems with this idea. First, in sequence labeling, some predictions are more wrong than others: we may miss only one tag out of fifty, or we may get all fifty wrong. We would like our learning algorithm to be sensitive to this difference. Second, the number of constraints is equal to the number of possible labelings, which is exponentially large in the length of the sequence.

⁵This model is also known as a **max-margin Markov network** (Taskar et al., 2003), emphasizing that the scoring function is constructed from a sum of components, which are Markov independent.

The first problem can be addressed by adjusting the constraint to require larger margins for more serious errors. Let $c(\mathbf{y}^{(i)}, \hat{\mathbf{y}}) \geq 0$ represent the **cost** of predicting label $\hat{\mathbf{y}}$ when the true label is $\mathbf{y}^{(i)}$. We can then generalize the margin constraint,

$$\forall \mathbf{y} \neq \mathbf{y}^{(i)}, \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) - \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) \geq c(\mathbf{y}^{(i)}, \mathbf{y}). \quad [6.78]$$

This cost-augmented margin constraint specializes to the constraint in Equation 6.77 if we choose the delta function $c(\mathbf{y}^{(i)}, \mathbf{y}) = \delta(\mathbf{y}^{(i)} \neq \mathbf{y})$. For sequence labeling, we can instead use a structured cost function, such as the **Hamming cost**,

$$c(\mathbf{y}^{(i)}, \mathbf{y}) = \sum_{m=1}^M \delta(y_m^{(i)} \neq y_m). \quad [6.79]$$

With this cost function, we require that the true labeling be separated from the alternatives by a margin that is proportional to the number of incorrect tags in each alternative labeling. Other cost functions are possible as well.

The second problem is that the number of constraints is exponential in the length of the sequence. This can be addressed by focusing on the prediction $\hat{\mathbf{y}}$ that *maximally* violates the margin constraint. We find this prediction by solving the following **cost-augmented decoding** problem:

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \neq \mathbf{y}^{(i)}} \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) - \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) + c(\mathbf{y}^{(i)}, \mathbf{y}) \quad [6.80]$$

$$= \operatorname{argmax}_{\mathbf{y} \neq \mathbf{y}^{(i)}} \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) + c(\mathbf{y}^{(i)}, \mathbf{y}), \quad [6.81]$$

where in the second line we drop the term $\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$, which is constant in \mathbf{y} .

We can now formulate the margin constraint for sequence labeling,

$$\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) - \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \left(\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) + c(\mathbf{y}^{(i)}, \mathbf{y}) \right) \geq 0. \quad [6.82]$$

If the score for $\theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$ is greater than the cost-augmented score for all alternatives, then the constraint will be met. Therefore we can maximize over the entire set $\mathcal{Y}(\mathbf{w})$, meaning that we can apply Viterbi directly.⁶

The name “cost-augmented decoding” is due to the fact that the objective includes the standard decoding problem, $\max_{\hat{\mathbf{y}}} \theta \cdot \mathbf{f}(\mathbf{w}, \hat{\mathbf{y}})$, plus an additional term for the cost. Essentially, we want to train against predictions that are strong and wrong: they should

⁶To maximize over the set $\mathcal{Y}(\mathbf{w}) \setminus \mathbf{y}^{(i)}$ we would need an alternative version of Viterbi that returns the k -best predictions. K -best Viterbi may be useful for other reasons — for example, in interactive applications, it can be helpful to show the user multiple possible taggings. The design of k -best Viterbi is left an exercise.

score highly according to the model, yet incur a large loss with respect to the ground truth. We can then adjust the weights to reduce the score of these predictions.

For cost-augmented decoding to be tractable, the cost function must decompose into local parts, just as the feature function $\mathbf{f}(\cdot)$ does. The Hamming cost, defined above, obeys this property. To solve this cost-augmented decoding problem using the Hamming cost, we can simply add features $f_m(y_m) = \delta(y_m \neq y_m^{(i)})$, and assign a weight of 1 to these features. Decoding can then be performed using the Viterbi algorithm. Are there cost functions that do not decompose into local parts? Suppose we want to assign a constant loss c to any prediction $\hat{\mathbf{y}}$ in which k or more predicted tags are incorrect, and zero loss otherwise. This loss function is combinatorial over the predictions, and thus we cannot decompose it into parts.

As with large-margin classifiers, it is possible to formulate the learning problem in an unconstrained form, by combining a regularization term on the weights and a Lagrangian for the constraints:

$$\min_{\boldsymbol{\theta}} \frac{1}{2} \|\boldsymbol{\theta}\|_2^2 - C \left(\sum_i \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}) - \max_{\mathbf{y} \in \mathcal{Y}(\mathbf{w}^{(i)})} [\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{y}) + c(\mathbf{y}^{(i)}, \mathbf{y})] \right), \quad [6.83]$$

In this formulation, C is a parameter that controls the tradeoff between the regularization term and the margin constraints. A number of optimization algorithms have been proposed for structured support vector machines, some of which are discussed in § 1.3.2. An empirical comparison by Kummerfeld et al. (2015) shows that stochastic subgradient descent — which is relatively easy to implement — is highly competitive, especially on the sequence labeling task of named entity recognition.

6.5.3 Conditional random fields

Structured perceptron is easy to implement, and structured support vector machines give excellent performance. However, sometimes we need to compute probabilities over labelings, $p(\mathbf{y} \mid \mathbf{w})$, and we would like to do this in a discriminative way. The **Conditional Random Field** (CRF; Lafferty et al., 2001) is a conditional probabilistic model for sequence labeling; just as structured perceptron is built on the perceptron classifier, conditional random fields are built on the logistic regression classifier.⁷ The basic probability model is,

$$p(\mathbf{y} \mid \mathbf{w}) = \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}))}{\sum_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}'))}. \quad [6.84]$$

⁷The name “Conditional Random Field” is derived from **Markov random fields**, a general class of models in which the probability of a configuration of variables is proportional to a product of scores across pairs (or more generally, cliques) of variables in a **factor graph**. In sequence labeling, the pairs of variables include all adjacent tags $\langle y_m, y_{m-1} \rangle$. The probability is **conditioned** on the words $\mathbf{w}_{1:M}$, which are always observed, motivating the term “conditional” in the name.

This is almost identical to logistic regression, but because the label space is now tag sequences, we require efficient algorithms for both **decoding** (searching for the best tag sequence given a sequence of words w and a model θ) and for **normalizing** (summing over all tag sequences). These algorithms will be based on the usual locality assumption on the feature function, $f(w, y) = \sum_{m=1}^M f(w, y_m, y_{m-1}, m)$.

Decoding in CRFs

Decoding — finding the tag sequence \hat{y} that maximizes $p(y | w)$ — is a direct application of the Viterbi algorithm. The key observation is that the decoding problem does not depend on the denominator of $p(y | w)$,

$$\begin{aligned} \hat{y} &= \operatorname{argmax}_y p(y | w) \\ &= \operatorname{argmax}_y \log p(y | w) \\ &= \operatorname{argmax}_y \theta \cdot f(y, w) - \log \sum_{y' \in \mathcal{Y}(w)} e^{\theta \cdot f(y', w)} \\ &= \operatorname{argmax}_y \theta \cdot f(y, w). \end{aligned}$$

This is identical to the decoding problem for structured perceptron, so the same Viterbi recurrence as defined in Equation 6.26 can be used.

Learning in CRFs

As with logistic regression, we learn the weights θ by minimizing the regularized negative log conditional probability,

$$\ell = \frac{\lambda}{2} \|\theta\|^2 - \sum_{i=1}^N \log p(y^{(i)} | w^{(i)}; \theta) \quad [6.85]$$

$$= \frac{\lambda}{2} \|\theta\|^2 - \sum_{i=1}^N \theta \cdot f(w^{(i)}, y^{(i)}) + \log \sum_{y' \in \mathcal{Y}(w^{(i)})} \exp(\theta \cdot f(w^{(i)}, y')), \quad [6.86]$$

where λ controls the amount of regularization. We will optimize θ by moving along the gradient of this loss. Probabilistic programming environments, such as THEANO (Bergstra et al., 2010) and TORCH (Collobert et al., 2011a), can compute this gradient using automatic differentiation. However, it is worth deriving the gradient to understand how this model works, and why learning is computationally tractable.

(c) Jacob Eisenstein 2014-2017. Work in progress.

As in logistic regression, the gradient includes a difference between observed and expected feature counts:

$$\frac{d\ell}{d\theta_j} = \lambda\theta_j + \sum_{i=1}^N E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] - f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)}), \quad [6.87]$$

where $f_j(\mathbf{w}^{(i)}, \mathbf{y}^{(i)})$ refers to the count of feature j for token sequence $\mathbf{w}^{(i)}$ and tag sequence $\mathbf{y}^{(i)}$.

The expected feature counts are computed by summing over all possible labelings of the word sequence,

$$E[f_j(\mathbf{w}^{(i)}, \mathbf{y})] = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w}^{(i)})} p(\mathbf{y} \mid \mathbf{w}^{(i)}; \boldsymbol{\theta}) f_j(\mathbf{w}^{(i)}, \mathbf{y}) \quad [6.88]$$

This looks bad: it is a sum over an exponential number of labelings. To solve this problem, we again rely on the assumption that the overall feature vector decomposes into a sum of local feature vectors, which we exploit to compute the expected feature counts as a sum across the sequence:

$$E[f_j(\mathbf{w}, \mathbf{y})] = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} p(\mathbf{y} \mid \mathbf{w}; \boldsymbol{\theta}) f_j(\mathbf{w}, \mathbf{y}) \quad [6.89]$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} p(\mathbf{y} \mid \mathbf{w}; \boldsymbol{\theta}) \sum_{m=1}^M f_j(\mathbf{w}, y_m, y_{m-1}, m) \quad [6.90]$$

$$= \sum_{m=1}^M \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} p(\mathbf{y} \mid \mathbf{w}; \boldsymbol{\theta}) f_j(\mathbf{w}, y_m, y_{m-1}, m) \quad [6.91]$$

$$= \sum_{m=1}^M \sum_{k, k' \in \mathcal{Y}} \sum_{\mathbf{y}: y_{m-1}=k', y_m=k} p(\mathbf{y} \mid \mathbf{w}; \boldsymbol{\theta}) f_j(\mathbf{w}, k, k', m) \quad [6.92]$$

$$= \sum_{m=1}^M \sum_{k, k' \in \mathcal{Y}} f_j(\mathbf{w}, k, k', m) \sum_{\mathbf{y}: y_{m-1}=k', y_m=k} p(\mathbf{y} \mid \mathbf{w}; \boldsymbol{\theta}) \quad [6.93]$$

$$= \sum_{m=1}^M \sum_{k, k' \in \mathcal{Y}} f_j(\mathbf{w}, k', k, m) \Pr(Y_{m-1} = k', Y_m = k \mid \mathbf{w}; \boldsymbol{\theta}). \quad [6.94]$$

This derivation works by interchanging the sum over tag sequences with the sum over indices m . At each position in the sequence, the locality restriction on features ensures that we need only the marginal probability of the tag bigram, $\Pr(Y_{m-1} = k', Y_m = k \mid \mathbf{w}; \boldsymbol{\theta})$. These tag bigram marginals are also used in unsupervised approaches to sequence labeling. In principle, these marginals still require a sum over the exponentially many label

sequences in which $Y_{m-1} = k'$ and $Y_m = k$. However, the marginals can be computed efficiently using the **forward-backward algorithm**.

*Forward-backward algorithm

Recall that in the hidden Markov model, it was possible to use the forward algorithm to compute marginal probabilities $p(y_m, \mathbf{w}_{1:m})$. We now derive a more general version of the forward algorithm, in which label sequences are scored in terms of **potentials** $\psi_m(k, k')$:

$$\psi_m(k, k') \triangleq \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, k, k', m)) \quad [6.95]$$

$$p(\mathbf{y} \mid \mathbf{w}) = \frac{\prod_{m=1}^M \psi_m(y_m, y_{m-1})}{\sum_{\mathbf{y}'} \prod_{m=1}^M \psi_m(y'_m, y'_{m-1})}. \quad [6.96]$$

Equation 6.96 simply expresses the CRF conditional likelihood, under a change of notation.

The tag bigram marginal probabilities can be written as,

$$\Pr(Y_{m-1} = k', Y_m = k \mid \mathbf{w}; \boldsymbol{\theta}) = \frac{\sum_{\mathbf{y}: y_m=k, y_{m-1}=k'} \prod_{n=1}^M \psi_n(y_n, y_{n-1})}{\sum_{\mathbf{y}'} \prod_{n=1}^M \psi_n(y'_n, y'_{n-1})}. \quad [6.97]$$

where the denominator is the marginal $p(\mathbf{w}_{1:M}) = \sum_{\mathbf{y}} p(\mathbf{w}, \mathbf{y}_{1:M})$, sometimes known as the **partition function**.⁸ Let us now consider how to compute each of these terms efficiently.

Computing the numerator In Equation 6.97, we sum over all tag sequences that include the transition $(Y_{m-1} = k') \rightarrow (Y_m = k)$. Because we are only interested in sequences that include this arc, we can decompose this sum into three parts: the sum over **prefixes** $\mathbf{y}_{1:m-1}$, the transition $(Y_{m-1} = k') \rightarrow (Y_m = k)$, and the sum over **suffixes** $\mathbf{y}_{m:M}$,

$$\begin{aligned} \sum_{\mathbf{y}: Y_m=k, Y_{m-1}=k'} \prod_{n=1}^M \psi_n(y_n, y_{n-1}) &= \sum_{\mathbf{y}_{1:m-1}: y_{m-1}=k'} \prod_{n=1}^{m-1} \psi_n(y_n, y_{n-1}) \\ &\quad \times \psi_m(k, k') \\ &\quad \times \sum_{\mathbf{y}_{m:M}: y_m=k} \prod_{n=m+1}^M \psi_n(y_n, y_{n-1}). \end{aligned} \quad [6.98]$$

The result is product of three terms: a score for getting to the position $(Y_{m-1} = k')$, a score for the transition from k' to k , and a score for finishing the sequence from $(Y_m = k)$.

⁸The terminology of “potentials” and “partition functions” comes from statistical mechanics (Bishop, 2006).

Let us define the first term as a **forward variable**,

$$\alpha_m(k) \triangleq \sum_{\mathbf{y}_{1:m}:y_m=k} \prod_{n=1}^m \psi_n(y_n, y_{n-1}) \quad [6.99]$$

$$= \sum_{k' \in \mathcal{Y}} \psi_m(k, k') \sum_{\mathbf{y}_{1:m-1}:y_{m-1}=k'} \prod_{n=1}^{m-1} \psi_n(y_n, y_{n-1}) \quad [6.100]$$

$$= \sum_{k' \in \mathcal{Y}} \psi_m(k, k') \times \alpha_{m-1}(k'). \quad [6.101]$$

Thus, we compute the forward variables while moving from left to right over the trellis. This forward recurrence is a generalization of the forward recurrence defined in § 6.4. If $\psi_m(k, k') = p_E(w_m | Y_m = k) \times p_T(k | k')$, then we exactly recover the Hidden Markov Model forward variable $\alpha_m(k) = p(\mathbf{w}_{1:m}, Y_m = k)$ as computed in § 6.4.3.

The third term of Equation 6.98 can also be defined recursively, this time moving over the trellis from right to left. The resulting recurrence is called the **backward algorithm**:

$$\beta_{m-1}(k) \triangleq \sum_{\mathbf{y}_{m-1:M}:y_{m-1}=k} \prod_{n=m}^M \psi_n(y_n, y_{n-1}) \quad [6.102]$$

$$= \sum_{k' \in \mathcal{Y}} \psi_m(k', k) \sum_{\mathbf{y}_{m:M}:y_m=k'} \prod_{n=m+1}^M \psi_n(y_n, y_{n-1}) \quad [6.103]$$

$$= \sum_{k' \in \mathcal{Y}} \psi_m(k', k) \times \beta_m(k'). \quad [6.104]$$

In practice, numerical stability requires that we use log-potentials rather than potentials, $\log \psi_m(y_m, y_{m-1}) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)$. Then the sums must be replaced with log-sum-exp:

$$\log \alpha_m(k) = \log \sum_{k' \in \mathcal{Y}} \exp(\log \psi_m(k, k') + \log \alpha_{m-1}(k')) \quad [6.105]$$

$$\log \beta_{m-1}(k) = \log \sum_{k' \in \mathcal{Y}} \exp(\log \psi_m(k', k) + \log \beta_m(k')). \quad [6.106]$$

Both the forward and backward algorithm operate on the trellis, which implies a space complexity $\mathcal{O}(MK)$. Because they require computing a sum over K terms at each node in the trellis, their time complexity is $\mathcal{O}(MK^2)$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Computing the normalization term The normalization term (partition function), sometimes abbreviated as Z , can be written as,

$$Z \triangleq \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} p(\mathbf{w}, \mathbf{y}) \quad [6.107]$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y})) \quad [6.108]$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \prod_{m=1}^M \exp(\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_m, y_{m-1}, m)) \quad [6.109]$$

$$= \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \prod_{m=1}^M \psi_m(y_m, y_{m-1}). \quad [6.110]$$

This term can be computed directly from either the forward or backward probabilities:

$$Z = \sum_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \prod_{m=1}^M \psi_m(y_m, y_{m-1}) \quad [6.111]$$

$$= \alpha_{M+1}(\diamond) \quad [6.112]$$

$$= \beta_0(\diamond). \quad [6.113]$$

CRF learning: wrapup Having computed the forward and backward variables, we can compute the desired marginal probability as,

$$P(Y_{m-1} = k', Y_m = k \mid \mathbf{w}_{1:M}) = \frac{\alpha_{m-1}(k') \psi_m(k, k') \beta_m(k)}{Z}. \quad [6.114]$$

This computation is known as the **forward-backward algorithm**. From the resulting marginals, we can compute the feature expectations $E[f_j(\mathbf{w}, \mathbf{y})]$; from these expectations, we compute a gradient on the weights $\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}}$. Stochastic gradient descent or quasi-Newton optimization can then be applied. As the optimization algorithm changes the weights, the potentials change, and therefore so do the marginals. Each iteration of the optimization algorithm therefore requires recomputing the forward and backward variables for each training instance.⁹

6.5.4 Neural sequence labeling

Recently, neural network methods have been applied to sequence labeling. These methods can be employed in tandem with structure prediction algorithms such as Viterbi,

⁹The `CRFsuite` package implements several learning algorithms for CRFs (<http://www.chokkan.org/software/crfsuite/>).

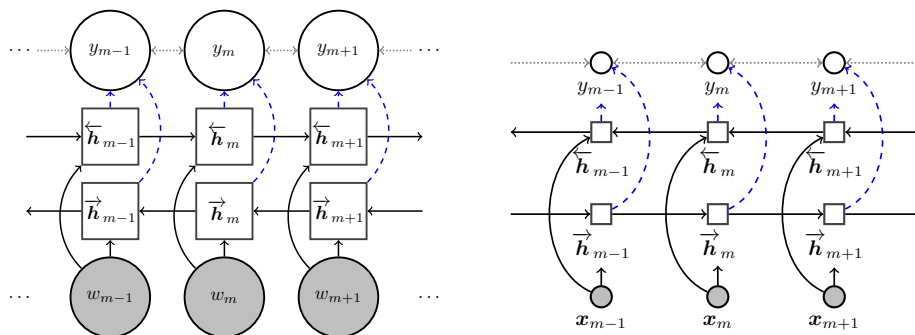


Figure 6.3: Bidirectional LSTM for sequence labeling. The solid lines indicate computation, the dashed lines indicate probabilistic dependency, and the dotted lines indicate the optional additional probabilistic dependencies between labels in the biLSTM-CRF.

although they are effective on their own. A relatively straightforward approach is to train a recurrent neural network or LSTM, as described in chapter 5; however, rather than predicting the next word in sequence, the model can be trained to predict the tag of the current word. A particularly effective approach is to train a **bidirectional recurrent neural network** (Graves and Schmidhuber, 2005), which estimates two hidden vectors, $\mathbf{h}_m = (\vec{\mathbf{h}}_m, \overleftarrow{\mathbf{h}}_m)$, with $\vec{\mathbf{h}}_m$ indicating the hidden state in a standard left-to-right RNN or LSTM, and $\overleftarrow{\mathbf{h}}_m$ indicating the hidden state in a left-to-right model. In this way, information from the entire sentence is brought to bear on the tagging decision for y_m . This approach was employed by Ling et al. (2015b), who find that bi-LSTMs perform considerably better than bi-RNNs on part-of-speech tagging, and that bidirectional variants of both models perform slightly but consistently better than their unidirectional counterparts. Note that LSTM-based tagging is a classification approach, so dynamic programming is not required to find the best tag sequence (this corresponds to Figure 6.3, excluding the finely dotted lines between adjacent tags).

It is also possible in neural sequence labeling to couple the tagging decisions, by introducing additional parameters for tag-to-tag transitions. This model is called the **LSTM-CRF**, due to its combination of aspects of the long short-term memory and conditional random field models (Huang et al., 2015). The model is shown in Figure 6.3. Lample et al. (2016) find that the LSTM-CRF is especially effective on the task of **named entity recognition**, a sequence labeling task that is described in detail in § 7.3. This task has particularly strong dependencies between adjacent tags, so it is not surprising to see an advantage for structure predictions here.

Both Ling et al. (2015b) and Lample et al. (2016) find it advantageous to model unseen and rare words through character-level representations. They train nested bi-LSTMs for each word, and take the concatenation of the introductory and final hidden states,

$(\vec{h}_M, \overleftarrow{h}_0)$ as the word embeddings, which is then used as input in the tagging bi-LSTM. Lample et al. (2016) combine these character-level embeddings with **pre-trained** word embeddings, which were estimated from an unlabeled dataset that is many orders of magnitude larger than the labeled data. They do not backpropagate into the word embeddings, and learn only the transition and output parameters of the model.

Convolutional Neural Networks for Sequence Labeling One disadvantage of recurrent neural networks is that the architecture requires iterating through the sequence of inputs and predictions: each hidden vector h_m must be computed from the previous hidden vector h_{m-1} , before predicting the tag y_m . These iterative computations are difficult to parallelize, and fail to exploit the speedups offered by **graphics processing units (GPUs)** on operations such as matrix multiplication. In contrast, **convolutional neural nets** predict each label y_m from a set of matrix operations on the neighboring word embeddings, $x_{m-k:m+k}$. Because there is no hidden state to update, the predictions for each y_m can be computed in parallel. For more on convolutional neural networks, see § 12.2.3. [todo: maybe move that section up here]

6.6 *Unsupervised sequence labeling

In unsupervised sequence labeling, we want to induce a Hidden Markov Model from a corpus of unannotated text $w^{(1)}, w^{(2)}, \dots, w^{(N)}$. 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 they might help us to better understand the language's inherent structure.

Unsupervised learning in hidden Markov models can be performed using the **Baum-Welch algorithm**, which combines forward-backward with expectation-maximization (EM). In the M-step, we compute the HMM parameters from expected counts:

$$\begin{aligned} \Pr(W = i \mid Y = k) &= \phi_{k,i} = \frac{E[\text{count}(W = i, Y = k)]}{E[\text{count}(Y = k)]} \\ \Pr(Y_m = k \mid Y_{m-1} = k') &= \lambda_{k',k} = \frac{E[\text{count}(Y_m = k, Y_{m-1} = k')]}{E[\text{count}(Y_{m-1} = k')]} \end{aligned}$$

The expected counts are computed in the E-step, using the forward and backward variables as defined in Equation 6.101 and Equation 6.104. Because we are working in a hidden Markov model, we define the potentials as,

$$\psi_m(k, k') = p_E(w_m \mid Y_m = k; \phi) \times p_T(Y_m = k \mid Y_{m-1} = k'; \lambda). \quad [6.115]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

The expected counts are then,

$$E[\text{count}(W = i, Y = k)] = \sum_{m=1}^M \Pr(Y_m = k \mid \mathbf{w}_{1:M}) \delta(W_m = i) \quad [6.116]$$

$$= \sum_{m=1}^M \frac{\Pr(Y_m = k, \mathbf{w}_{1:m}) p(\mathbf{w}_{m+1:M} \mid Y_m = k)}{p(\mathbf{w}_{1:M})} \delta(W_m = i) \quad [6.117]$$

$$= \frac{1}{\alpha_{M+1}(\diamond)} \sum_{m=1}^M \alpha_m(k) \beta_m(k) \delta(W_m = i) \quad [6.118]$$

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 $p(\mathbf{w}_{1:M}) = \alpha_{M+1}(\diamond) = \beta_0(\diamond)$.

The expected transition counts can be computed in a similar manner:

$$E[\text{count}(Y_m = k, Y_{m-1} = k')] = \sum_{m=1}^M \Pr(Y_m = k, Y_{m-1} = k' \mid \mathbf{w}_{1:M}) \quad [6.119]$$

$$\propto \sum_{m=1}^M p(Y_{m-1} = k', \mathbf{w}_{1:m-1}) p(w_{m+1:M} \mid Y_m = k) \\ \times p(w_m, Y_m = k \mid Y_{m-1} = k') \quad [6.120]$$

$$= \sum_{m=1}^M p(Y_{m-1} = k', \mathbf{w}_{1:m-1}) p(w_{m+1:M} \mid Y_m = k) \\ \times p(w_m \mid Y_m = k) p(Y_m = k \mid Y_{m-1} = k') \quad [6.121]$$

$$= \sum_{m=1}^M \alpha_{m-1}(k') \beta_m(k) \phi_{k,w_m} \lambda_{k' \rightarrow k}. \quad [6.122]$$

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 $\Pr(Y_{m-1} = k', Y_m = k \mid \mathbf{w}_{1:M})$.

6.6.1 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**

(c) Jacob Eisenstein 2014-2017. Work in progress.

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 trellis of values at each step, 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).

6.6.2 Alternative unsupervised learning methods

As noted in § 4.5, 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):

$$p(y_m \mid \mathbf{y}_{-m}, \mathbf{w}_{1:M}) \propto p(w_m \mid y_m)p(y_m \mid \mathbf{y}_{-m}). \quad [6.123]$$

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).

Exercises

1. Consider the garden path sentence, *The old man the boat*. Given word-tag and tag-tag features, what inequality in the weights must hold for the correct tag sequence to outscore the garden path tag sequence for this example?
2. Sketch out an algorithm for a variant of Viterbi that returns the top- n label sequences. What is the time and space complexity of this algorithm?
3. Show how to compute the marginal probability $\Pr(Y_{m-2} = k, Y_m = k')$, in terms of the forwards and backward variables, and the potentials ψ_m .

(c) Jacob Eisenstein 2014-2017. Work in progress.

4. more tk

Chapter 7

Applications of sequence labeling

Sequence labeling has applications throughout natural language processing. This chapter focuses on the classical applications of part-of-speech tagging, morpho-syntactic attribute tagging, named entity recognition, and tokenization. It also touches briefly on two applications to interactive settings: dialogue act recognition and the detection of code-switching points between languages.

7.1 Part-of-speech tagging

The **syntax** of a language is the collection of principles under which sequences of words are judged to be grammatically acceptable by fluent speakers. One of the most basic syntactic concepts is the **part-of-speech** (POS), which refers to the syntactic role of each word in a sentence. We have already referred to this concept informally in the previous chapter, and you may have some intuitions from your own study of English. For example, in the sentence *We like vegetarian sandwiches*, you may already know that *we* and *sandwiches* are nouns, *like* is a verb, and *vegetarian* is an adjective. These labels depend on the context in which the word appears: in *she eats like a vegetarian*, the word *like* is a preposition, and the word *vegetarian* is a noun.

Parts-of-speech can help to disentangle or explain various linguistic problems. Recall Chomsky's proposed distinction in chapter 5:

(7.1) Colorless green ideas sleep furiously.

(7.2) *Ideas colorless furiously green sleep.

Why is the first grammatical and the second not? One explanation is that the first example contains part-of-speech transitions that are typical in English: adjective to adjective, adjective to noun, noun to verb, and verb to adverb. In contrast, the second sentence contains transitions that are unusual: noun to adjective and adjective to verb. The ambiguity in

sentences like *teacher strikes idle children* can also be explained in terms of parts of speech: in the interpretation that was likely intended, *strikes* is a noun and *idle* is a verb; in the alternative explanation, *strikes* is a verb and *idle* is an adjective.

Part-of-speech tagging is often taken as a early step in a natural language processing pipeline. Indeed, parts-of-speech provide features that can be useful for many of the tasks that we will encounter later, such as parsing, coreference resolution, and relation extraction.

7.1.1 Parts-of-Speech

The **Universal Dependencies** project (UD) is an effort to create syntactically-annotated corpora across many languages, using a single annotation standard (Nivre et al., 2016). As part of this effort, they have designed a part-of-speech tag inventory, which is meant to capture word classes across as many languages as possible.¹ This section describes that inventory, giving rough definitions for each of tags, along with supporting examples.

Part-of-speech tags are **morphosyntactic**, rather than **semantic**, categories. This means that they describe words in terms of how they pattern together and how they are internally constructed (e.g., what suffixes and prefixes they include). For example, you may think of a noun as referring to objects or concepts, and verbs as referring to actions or events. But events can also be nouns — for example, *the **running** of the bulls*, where the determiner *the* is a strong clue that *running* is acting as a noun.

The Universal Dependency part-of-speech tagset

The UD tagset is broken up into three groups: open class tags, closed class tags, and “other”.

Open class tags Nearly all languages contain nouns, verbs, adjectives, and adverbs.² These are all **open word classes**, because new words are constantly being added to them. The UD tagset includes two other tags that are open classes: proper nouns and interjections.

- **Nouns** (UD tag: NOUN) tend to describe entities and concepts, e.g.,

(7.3) **Toes** are scarce among veteran **blubber men**.

In English, nouns tend to follow determiners and adjectives, and can play the subject role in the sentence. They can be marked for the plural number by an -s suffix.

¹The UD tagset builds on earlier work from Petrov et al. (2012), in which a set of twelve universal tags was identified by creating mappings from tagsets for individual languages.

²One well-known exception is Korean, which some linguists argue does not have adjectives Kim (2002).

- **Proper nouns** (PROPN) are tokens in names, which uniquely specify a given entity,

(7.4) “**Moby Dick?**” shouted **Ahab**.

- **Verbs** (VERB) describe actions or events

(7.5) “Moby Dick?” **shouted** Ahab.

(7.6) **Shall** we **keep chasing** this murderous fish?

English verbs tend to come in between the subject and some number of direct objects, depending on the verb. They can be marked for **tense** and **aspect** using suffixes such as *-ed* and *-ing*. (These suffixes are an example of **inflectional morphology**, which is discussed in more detail in § 8.1.4.)

- **Adjectives** (ADJ) describe properties of entities,

(7.7) Shall we keep chasing this **murderous** fish?

(7.8) Toes are **scarce** among **veteran** blubber men.

In the second example, *scarce* is a predicative adjective, linked to the subject by the **copula verb** *are*. In contrast, *murderous* and *veteran* are attribute adjectives, modifying the noun phrase in which they are embedded.

- **Adverbs** (ADV) describe properties of events, and may also modify adjectives or other adverbs:

(7.9) It is not down on any map; true places **never** are.

(7.10) ... **treacherously** hidden beneath the loveliest tints of azure

(7.11) Not drowned **entirely**, though.

- **Interjections** (INTJ) are used in exclamations, e.g.,

(7.12) **Aye aye!** it was that accursed white whale that razed me.

Closed class tags Closed word classes rarely receive new members. They are sometimes referred to as **function words** — as opposed to **content words** — as they have little lexical meaning of their own, but rather, help to organize the components of the sentence.

- **Adpositions** (ADP) describe the relationship between a complement (usually a noun phrase) and another unit in the sentence, typically a noun or verb phrase.

(7.13) Toes are scarce **among** veteran blubber men.

(c) Jacob Eisenstein 2014-2017. Work in progress.

(7.14) It is not **down on** any map.

(7.15) Give not thyself **up** then.

As the examples show, English generally uses prepositions, which are adpositions that appear before their complement. (An exception is *ago*, as in, *we met three days ago*). Postpositions are used in other languages, such as Japanese and Turkish.

- **Auxiliary verbs** (AUX) are a closed class of verbs that add information such as tense, aspect, person, and number.

(7.16) **Shall** we keep chasing this murderous fish?

(7.17) What the white whale was to Ahab, **has been** hinted.

(7.18) Ahab **must** use tools

(7.19) Meditation and water **are** wedded forever.

(7.20) Toes **are** scarce among veteran blubber men.

As noted above, the final example is a copula verb, which is also tagged as an auxiliary in the UD corpus.

- **Coordinating conjunctions** (CCONJ) express relationships between two words or phrases, which play a parallel role:

(7.21) Meditation **and** water are wedded forever.

- **Subordinating conjunctions** (SCONJ) link two elements, making one syntactically subordinate to the other:

(7.22) There is wisdom **that** is woe.

- **Pronouns** (PRON) are words that substitute for nouns or noun phrases when the meaning is clear from context.

(7.23) Be **it what it** will, I 'll go to **it** laughing.

(7.24) I try all things, I achieve **what I** can.

The example includes the personal pronouns *I* and *it*, as well as the relative pronoun *what*. Other pronouns include *myself*, *somebody*, and *nothing*.

- **Determiners** (DET) provide additional information about the nouns or noun phrases that they modify:

(7.25) What **the** white whale was to Ahab, has been hinted.

(7.26) It is not down on **any** map.

(c) Jacob Eisenstein 2014-2017. Work in progress.

(7.27) I try **all** things ...

(7.28) Shall we keep chasing **this** murderous fish?

Determiners include articles (*the*), possessive determiners (*their*), demonstratives (*this murderous fish*), and quantifiers (*any map*).

- **Numerals** (NUM) are an infinite but closed class, which includes integers, fractions, and decimals, regardless of whether spelled out or written in numerical form.

(7.29) How then can this **one** small heart beat

(7.30) I am going to put him down for the **three hundredth**

- **Particles** (PART) are a catch-all of function words that combine with other words or phrases, but do not meet the conditions of the other tags. In English, this includes the infinitival *to*, the possessive marker, and negation.

(7.31) Better **to** sleep with a sober cannibal than a drunk Christian.

(7.32) So man **'s** insanity is heaven **'s** sense

(7.33) It is **not** down on any map

As the first example shows, the possessive marker is not considered part of the same token as the word that it modifies, so that *man's* is split into two tokens. (Tokenization is described in more detail in § 7.4.) Non-English examples of particles include question markers, such as

In other tagsets, particularly those designed for English like the Penn Treebank, verbal particles such as *give up* and *fuck off* are tagged as particles. However, in the UD corpus, these are tagged as adpositions.

Other Remaining UD tags include punctuation (PUN) and symbols (SYM). Punctuation is purely structural — e.g., commas, periods, colons — while symbols can carry content of their own. Examples of symbols include dollar and percentage symbols, mathematical operators, emoticons, emojis, and internet addresses. A final catch-all tag is X, which is used for words that cannot be assigned another part-of-speech category. The X tag is also used in cases of **code-switching** (between languages), described in § 7.5.

Other tagsets

Prior to the Universal Dependency treebank, part-of-speech tagging was performed using language-specific tagsets. The dominant tagset for English is the Penn Treebank (PTB). The PTB tagset includes 45 tags — much more granularity than the UD tagset. This granularity is reflected in distinctions between singular and plural nouns, various verb tenses

and aspects, possessive and non-possessive pronouns, comparative and superlative adjectives and adverbs (e.g., *faster*, *fastest*), and so on. The Brown corpus includes a tagset that is even more detailed, with 87 tags Francis (1964), including special tags for individual auxiliary verbs such as *be*, *do*, and *have*.

Different languages make different distinctions, and so the PTB and Brown tagsets are not appropriate for a language such as Chinese, which does not mark the verb tense (Xia, 2000); nor for Spanish, which marks every combination of person and number in the verb ending; nor for German, which marks the case of each noun phrase. Each of these languages requires more detail than English in some areas of the tagset, and less in other areas. The strategy of the Universal Dependency corpus is to design a coarse-grained tagset to be used across all languages, and then to additionally annotate language-specific **morphosyntactic attributes**, such as number, tense, and case. The attribute tagging task is described in more detail in § 7.2.

Finally, social media such as Twitter have been shown to require tagsets of their own (Gimpel et al., 2011). Such corpora contain some tokens that are simply not equivalent to anything encountered in a typical written corpus: e.g., emoticons, URLs, and hashtags. Social media also includes dialectal words like *gonna* (going to, e.g. *We gonna be fine*) and *Ima* (I'm going to, e.g., *Ima tell you one more time*), which can be analyzed either as non-standard orthography (making tokenization impossible), or as lexical items in their own right. In either case, it is clear that existing tags like NOUN and VERB cannot handle cases like *Ima*, which combine aspects of the noun and verb. Gimpel et al. (2011) therefore propose a new set of tags to deal with these cases.

7.1.2 Accurate part-of-speech tagging

Part-of-speech tagging is the problem of selecting the correct tag for each word in a sentence. Success is typically measured by accuracy on an annotated test set, which is simply the fraction of tokens that were tagged correctly.

Baselines

A simple baseline for part-of-speech tagging is to choose the most common tag for each word. For example, in the Universal Dependency treebank, the word *talk* appears 96 times, and 85 of those times it is labeled as a verb: therefore, this baseline will always predict verb when we see it in the test set. For words that do not appear in the training corpus, the baseline simply guesses the most common tag overall, which is noun. In the Penn Treebank, this simple baseline obtains accuracy above 92%. A more rigorous evaluation is the accuracy on **out-of-vocabulary words**, which are not seen in the training data. Tagging the words correctly requires attention to the context or the word's internal structure.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Contemporary approaches

Conditional random fields and structured perceptron perform at or near the state-of-the-art for part-of-speech tagging in English. For example, (Collins, 2002) achieved 97.1% accuracy on the Penn Treebank, using a structured perceptron, using the following base features (originally introduced by Ratnaparkhi (1996)):

- current word, w_m
- previous words, w_{m-1}, w_{m-2}
- next words, w_{m+1}, w_{m+2}
- previous tag, y_{m-1}
- previous two tags, (y_{m-1}, y_{m-2})
- for rare words:
 - first k characters, up to $k = 4$
 - last k characters, up to $k = 4$
 - whether w_m contains a number, uppercase character, or hyphen.

Similar results for the PTB data have been achieved using conditional random fields (CRFs; Toutanova et al., 2003).

More recent work has demonstrated the power of neural sequence models, such as the **long short-term memory (LSTM)** (§ 6.5.4). Plank et al. (2016) apply a CRF and a bidirectional LSTM to twenty-two languages in the UD corpus, achieving an average accuracy of 94.3% for the CRF, and 96.5% with the bi-LSTM. Their model employs three types of embeddings: learned word embeddings, which are updated during training; pre-trained word embeddings, which are never updated, but which help to tag out-of-vocabulary words; and character-based embeddings. These character-based embeddings are computed by running an LSTM on the individual characters in each word, thereby capturing common orthographic patterns such as prefixes, suffixes, and capitalization. Extensive evaluations show that these additional embeddings are crucial to their model’s success.

7.2 Morphosyntactic Attributes

The Universal Dependency (UD) tagset provides a coarse-grained description of the syntactic role of each word in a sentence, using a set of tags that generalizes across languages. However, there is considerably more to say about a word than whether it is a noun or a verb: in English, verbs are distinguished by features such as tense and aspect, nouns by number, adjectives by degree, and so on. Unfortunately, these features are language-specific: other languages distinguish other features, such as **case** (the role of the noun with respect to the action of the sentence, which is marked in many languages, such as Latin and Ger-

(c) Jacob Eisenstein 2014-2017. Work in progress.

word	PTB tag	UD tag	UD attributes
<i>The</i>	DT	DET	DEFINITE=DEF PRONTYPE=ART
<i>German</i>	JJ	ADJ	DEGREE=POS
<i>Expressionist</i>	NN	NOUN	NUMBER=SING
<i>movement</i>	NN	NOUN	NUMBER=SING
<i>was</i>	VBD	AUX	MOOD=IND NUMBER=SING PERSON=3 TENSE=PAST VERBFORM=FIN
<i>destroyed</i>	VBN	VERB	TENSE=PAST VERBFORM=PART VOICE=PASS
<i>as</i>	IN	ADP	
<i>a</i>	DT	DET	DEFINITE=IND PRONTYPE=ART
<i>result</i>	NN	NOUN	NUMBER=SING
<i>.</i>	.	PUNCT	

Figure 7.1: UD and PTB part-of-speech tags, and UD morphosyntactic attributes. Example selected from the UD 1.4 English corpus.

man³ and **evidentiality** (the source of information for the speaker’s statement, which is marked in languages such as Turkish). In the UD corpus, these attributes are annotated as feature-value pairs for each token.⁴

An example is shown in Figure 7.1. The determiner *the* is marked with two attributes: PRONTYPE=ART, which indicates that it is an **article** (as opposed to another type of determiner or pronominal modifier), and DEFINITE=DEF, which indicates that it is a **definite article** (referring to a specific, known entity). The verbs are each marked with several attributes. The auxiliary verb *was* is third-person, singular, past tense, finite (conjugated), and indicative (describing an event that has happened or is currently happenings); the main verb *destroyed* is in participle form (so there is no additional person and number information), past tense, and passive voice. Some, but not all, of these distinctions are reflected in the PTB tags VBD (past-tense verb) and VBN past participle.

³Case is marked in English for some personal pronouns, e.g., *She saw her*, *They saw them*.

⁴The annotation and tagging of morphosyntactic attributes can be traced back to earlier work on Turkish (Oflaz and Kuruöz, 1994) and Czech (Hajič and Hladká, 1998). MULTTEXT-East was an early multilingual corpus to include morphosyntactic attributes (Dimitrova et al., 1998).

While there are thousands of papers on part-of-speech tagging, there is comparatively little work on automatically labeling morphosyntactic attributes. Faruqui et al. (2016a) train a support vector machine classification model, using a minimal feature set that includes the word itself, its prefixes and suffixes, and type-level information listing all possible morphosyntactic attributes for each word and its neighbors. Mueller et al. (2013) use a conditional random field (CRF), in which the tag space consists of all observed combinations of morphosyntactic attributes (e.g., the tag would be DEF+ART for the word *the* in Figure 7.1). This massive tag space is managed by decomposing the feature space over individual attributes, and pruning paths through the trellis. More recent work has employed bidirectional LSTM sequence models. For example, Pinter et al. (2017) train a bidirectional LSTM sequence model. The input layer and hidden vectors in the LSTM are shared across attributes, but each attribute has its own output layer, culminating in a softmax over all attribute values, e.g. $y_t^{\text{NUMBER}} \in \{\text{SING}, \text{PLURAL}, \dots\}$. They find that character-level information is crucial, especially when the amount of labeled data is limited.

Evaluation is performed by first computing recall and precision for each attribute. These scores can then be averaged at either the type or token level to obtain micro- or macro- F -measure. Pinter et al. (2017) evaluate on 23 languages in the UD treebank, reporting a median micro- F -measure of 0.95. Performance is strongly correlated with the size of the labeled dataset for each language, with a few outliers: for example, Chinese is particularly difficult, because although the dataset is relatively large (10^5 tokens in UD 1.4), only 6% of tokens have any attributes, offering few useful labeled instances.

7.3 Named entity recognition

A core problem in information extraction is to recognize and extract mentions of **named entities** in text. In news documents, the core entity types are people, locations, and organizations; more recently, the task has been extended to include amounts of money, percentages, dates, and times. In (7.34) in Figure 7.2, the named entities include: *The U.S. Army*, an organization; *Atlanta*, a location; and *May 14, 1864*, a date. Named entity recognition is also a key task in **biomedical natural language processing**, with entity types including proteins, DNA, RNA, and cell lines (e.g., Collier et al., 2000; Ohta et al., 2002). Figure 7.2 shows an example from the GENIA corpus.

A standard approach to tagging named entity spans is to use discriminative sequence labeling methods such as conditional random fields and structured perceptron. As described in chapter 6, these methods use the Viterbi algorithm to search over all possible label sequences, while scoring each sequence using a feature function that decomposes across adjacent tags. Named entity recognition is formulated as a tagging problem by assigning each word token to a tag from a tagset. However, there is a major difference from part-of-speech tagging: in NER we need to recover **spans** of tokens, such as *The*

(c) Jacob Eisenstein 2014-2017. Work in progress.

- (7.34) *The U.S. Army captured Atlanta on May 14 , 1864*
 B-ORG I-ORG I-ORG O B-LOC O B-DATE I-DATE I-DATE I-DATE
- (7.35) *Number of glucocorticoid receptors in lymphocytes and ...*
 O O B-PROTEIN I-PROTEIN O B-CELLTYPE O ...

Figure 7.2: BIO notation for named entity recognition. Example (7.35) is drawn from the GENIA corpus of biomedical documents (Ohta et al., 2002).

United States Army. To do this, the tagset must distinguish tokens that are at the **beginning** of a span from tokens that are **inside** a span.

This is accomplished by the **BIO notation**, shown in Figure 7.2. Each token at the beginning of a name span is labeled with a B- prefix; each token within a name span is labeled with an I- prefix. Tokens that are not parts of name spans are labeled as O. From this representation, it is unambiguous to recover the entity name spans within a labeled text. Another advantage is from the perspective of learning: tokens at the beginning of name spans may have different properties than tokens within the name, and the learner can exploit this. This insight can be taken even further, with special labels for the **last** tokens of a name span, and for **unique** tokens in name spans, such as *Atlanta* in the example in Figure 7.2. This is called **BILOU** notation, and has been shown to yield improvements in supervised named entity recognition (Ratinov and Roth, 2009).

Feature-based sequence labeling Named entity recognition was an early application of conditional random fields (McCallum and Li, 2003). The use of Viterbi decoding restricts the feature function $f(w, y) = \sum_m f(w, y_m, y_{m-1}, m)$, so that each feature can consider only local adjacent tags. Typical features include tag transitions, word features for w_m and its neighbors, character-level features for prefixes and suffixes, and “word shape” features to capture capitalization. As an example, base features for the word *Army* in the example in (7.34) include:

$\langle \text{CURR-WORD:} \textit{Army}, \text{PREV-WORD:} \textit{U.S.}, \text{NEXT-WORD:} \textit{captured}, \text{PREFIX-1:} \textit{A-},$
 $\text{PREFIX-2:} \textit{Ar-}, \text{SUFFIX-1:} \textit{-y}, \text{SUFFIX-2:} \textit{-my}, \text{SHAPE:} \textit{Xxxx} \rangle$

Another source of features is to use **gazetteers**: lists of known entity names. For example, the U.S. Social Security Administration provides a list of tens of thousands of given names — more than could be observed in any reasonable annotated corpus. Tokens or spans that match an entry in a gazetteer can receive special features; this provides a way to incorporate hand-crafted resources such as name lists in a learning-driven framework.

Neural sequence labeling for NER Current research has emphasized neural sequence labeling, using similar LSTM models to those employed in part-of-speech tagging (Ham-

[todo: import zh-example successfully]

Figure 7.3: Example of Chinese tokenization (Sproat et al., 1996)

merton, 2003; Huang et al., 2015; Lample et al., 2016). The bidirectional LSTM-CRF (Figure 6.3 in § 6.5.4) does particularly well on this task, due to its ability to model tag-to-tag dependencies. However, Strubell et al. (2017) show that **convolutional neural networks** can be equally accurate, with significant improvement in speed due to the efficiency of implementing ConvNets on **graphics processing units (GPUs)**. The key innovation in this work was the use of **dilated convolutions**, which are described in more detail in § 2.4.

7.4 Tokenization

A basic problem for text analysis, first discussed in § 3.3.1, is to break the text into a sequence of discrete tokens. For alphabetic languages such as English, deterministic scripts suffice to achieve accurate tokenization. However, in logographic writing systems such as Chinese script, words are typically composed of a small number of characters, without space in between (Figure 7.3). One approach is to match character sequences against a known dictionary (e.g., Sproat et al., 1996), using additional statistical information about word frequency. However, no dictionary is completely comprehensive, and dictionary-based approaches can struggle with such out-of-vocabulary words.

Chinese tokenization has therefore been approached as a supervised sequence labeling problem. Xue et al. (2003) train a logistic regression classifier to make independent segmentation decisions while moving a sliding window across the document. A set of rules is then used to convert these individual classification decisions into an overall tokenization of the input. However, these individual decisions may be globally suboptimal, motivating a structure prediction approach. Peng et al. (2004) train a conditional random field model to predict labels of START or NONSTART on each character. More recent work has employed neural network architectures. For example, Chen et al. (2015) use an LSTM-CRF architecture, as described in § 6.5.4: they construct a trellis, in which each tag is scored according to the hidden state of an LSTM, and tag-tag transitions are scored according to learned transition weights. The best-scoring segmentation is then computed by the Viterbi algorithm.

7.5 Code switching

[todo: needs work] Multilingual speakers and writers do not restrict themselves to a single language; rather, they may switch between languages, in a phenomenon known as **code switching** (Auer, 2013; Poplack, 1980). Code switching is particularly common in

(c) Jacob Eisenstein 2014-2017. Work in progress.

Speaker	Dialogue Act	Utterance
A	YES-NO-QUESTION	<i>So do you go college right now?</i>
A	ABANDONED	<i>Are yo-</i>
B	YES-ANSWER	<i>Yeah,</i>
B	STATEMENT	<i>It's my last year [laughter].</i>
A	DECLARATIVE-QUESTION	<i>You're a, so you're a senior now.</i>
B	YES-ANSWER	<i>Yeah,</i>
B	STATEMENT	<i>I'm working on my projects trying to graduate [laughter]</i>
A	APPRECIATION	<i>Oh, good for you.</i>
B	BACKCHANNEL	<i>Yeah.</i>

Figure 7.4: Example of dialogue act labeling (Stolcke et al., 2000).

An example is shown in Figure 7.4. The annotation is performed over **UTTERANCES**, with the possibility of multiple segments per **conversational turn** (in rare cases such as interruptions, an utterance may split over multiple turns). Some segments are clauses (e.g., *So do you go to college right now?*), while others are single words (e.g., *yeah*). Stolcke et al. (2000) report that hidden Markov models (HMMs) achieve 96% accuracy on supervised utterance segmentation. The labels themselves reflect the conversational goals of the speaker: the utterance *yeah* functions as an answer in response to the question *you're a senior now*, but as a **backchannel** (demonstrating comprehension and encouraging the other speaker to continue talking) in the final line of the excerpt.

For task of dialogue act labeling, Stolcke et al. (2000) apply a hidden Markov model. The generative probability $p(\mathbf{w}_m \mid \mathbf{y}_m)$ must generate the entire sequence of words in the utterance, and it is modeled as a trigram language model (§ 5.1). Stolcke et al. (2000) also account for acoustic features, which capture the **prosody** of each utterance — for example, tonal and rhythmic properties of speech, which can be used to distinguish dialogue acts such as questions and answers. These features are handled with an additional emission distribution, $p(\mathbf{a}_m \mid \mathbf{y}_m)$, which is modeled with a probabilistic decision tree (Murphy, 2012). They find that while acoustic features yield small improvements overall, they play an important role in distinguish questions from statements, and agreements from backchannels.

Recurrent neural architectures for dialogue act labeling have been proposed by Kalchbrenner and Blunsom (2013) and Ji et al. (2016), with strong empirical results. Both models are recurrent at the utterance level, so that each complete utterance updates a hidden state. The recurrent-convolutional network of Kalchbrenner and Blunsom (2013) uses convolution to obtain a representation of each individual utterance, while Ji et al. (2016) use a second level of recurrence, over individual words. This enables their method to also func-

forums.

(c) Jacob Eisenstein 2014-2017. Work in progress.

tion as a language model, giving probabilities over sequences of words in a document.

Chapter 8

Formal language theory

We have now seen methods for learning to label individual words, vectors of word counts, and sequences of words; we will soon proceed to more complex structural transformations. Most of these techniques could apply to counts or sequences from any discrete vocabulary; there is nothing fundamentally linguistic about, say, a hidden Markov model. This raises an embarrassingly basic question that this text has not yet considered: what is a language?

This chapter will take the perspective of **formal language theory**, in which a **language** is defined as a set of **strings**, each of which is a sequence of elements from a finite alphabet. For interesting languages, there are an infinite number of strings that are in the language, and an infinite number of strings that are not. For example:

- the set of all even-length sequences from the alphabet $\{a, b\}$, e.g., $\{\emptyset, aa, ab, ba, bb, aaaa, aaab, \dots\}$;
- the set of all sequences from the alphabet $\{a, b\}$ that contain *aaa* as a substring, e.g., $\{aaa, aaaa, baaa, aaab, \dots\}$;
- the set of all sequences of English words (drawn from a finite dictionary) that contain at least one verb (a finite subset of the dictionary);
- the `python` programming language.

Formal language theory has defined various classes of languages and their computational properties. In particular, theorists have delineated classes of languages based on the computational complexity of solving the **membership problem** — determining whether a string is in a language. The chapter will focus on three classes of formal languages: regular, context-free, and “mildly” context-sensitive languages.

A key insight of 20th century linguistics is that formal language theory can be usefully applied to natural languages such as English, by designing formal languages that capture as many properties of the natural language as possible. For many such formalisms, a

useful linguistic analysis comes as a byproduct of solving the membership problem. Furthermore, the membership problem can be generalized to the problems of **scoring** strings for their acceptability (as in language modeling), and of **transducing** one string into another.

8.1 Regular languages

Sooner or later, most computer scientists will write a **regular expression**. If you have, then you have defined a **regular language**, which is any language that can be defined by a regular expression. Formally, a regular expression can include the following elements:

- A **literal character** drawn from some finite alphabet Σ .
- The **empty string** ϵ .
- The concatenation of two regular expressions RS , where R and S are both regular expressions. The resulting expression accepts a string that can be decomposed $x = yz$, where y is accepted by R and z is accepted by S .
- The alternation $R \mid S$, where R and S are both regular expressions. The resulting expression accepts a string x if it is accepted by R or it is accepted by S .
- The **Kleene star** R^* , which accepts any string x that can be decomposed into a sequence of strings which are all accepted by R .
- Parenthesization (R) , which is used to limit the scope of the concatenation, alternation, and Kleene star operators.

Here are some example regular expressions for some of the languages described above:

- The set of all even length strings on the alphabet $\{a, b\}$: $((aa)|(ab)|(ba)|(bb))^*$
- The set of all sequences of the alphabet $\{a, b\}$ that contain aaa as a substring: $(a|b)^*aaa(a|b)^*$
- The set of all sequences of English words that contain at least one verb: W^*VW^* , where W is an alternation between all words in the dictionary, and V is an alternation between all verbs ($V \subseteq W$).

We do not include a regular expression for the Python programming language, because this language is not regular — there is no regular expression that can capture its syntax. We will discuss why towards the end of this section.

Regular languages are **closed** under union, intersection, and concatenation. This means, for example, that if two languages L_1 and L_2 are regular, then so are the languages $L_1 \cup L_2$, $L_1 \cap L_2$, and the language of strings that can be decomposed as $s = tu$, with $s \in L_1$ and $t \in L_2$. Regular languages are also closed under negation: if L is regular, then so is the language $\bar{L} = \{s \notin L\}$.

8.1.1 Finite-state acceptors

A regular expression defines a regular language, but does not give an algorithm for determining whether a string is in the language that it defines. **Finite-state automata** are theoretical models of algorithms for regular languages, which involve transitions between a finite number of states. The most basic type of finite-state automaton is the **finite-state acceptor (FSA)**, which describes the computation involved in testing if a string is a member of a language. Formally, a finite-state acceptor is 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** $Q = \{q_0, q_1, \dots, q_n\}$;
- a **start state** $q_0 \in Q$;
- a set of **final states** $F \subseteq Q$;
- a **transition function** $\delta : Q \times (\Sigma \cup \{\epsilon\}) \rightarrow 2^Q$. The transition function maps from a state and an input symbol (or empty string ϵ) to a **set** of possible resulting states.

A **path** in M is a sequence of transitions, $\pi = t_1, t_2, \dots, t_N$, where each t_i traverses an arc in the transition function δ . The finite-state acceptor M accepts a string ω if there is a **accepting path**, in which the initial transition t_1 begins at the start state q_0 , the final transition t_N terminates in a final state in Q , and the entire input ω is consumed.

Example

Consider the following FSA, M_1 .

$$\Sigma = \{a, b\} \quad [8.1]$$

$$Q = \{q_0, q_1\} \quad [8.2]$$

$$F = \{q_1\} \quad [8.3]$$

$$\begin{aligned} \delta = \{ & \{(q_0, a) \rightarrow q_0\}, \\ & \{(q_0, b) \rightarrow q_1\}, \\ & \{(q_1, b) \rightarrow q_1\} \} \end{aligned} \quad [8.4]$$

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 the machine is in state q_0 and reads symbol a , it stays in q_0 . A graphical state diagram representation for M_1 is shown in Figure 8.1. Because each pair of initial state and symbol has at most one resulting state, M_1 is **deterministic**: each string ω induces at most one accepting path. Note that there are no transitions for the symbol a in state q_1 ; if a is encountered in q_1 , then the acceptor is stuck, and the input string is rejected.

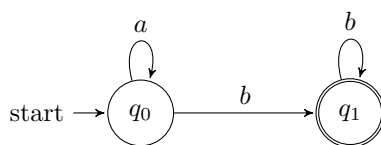


Figure 8.1: State diagram for the finite-state acceptor M_1 .

What strings does M_1 accept? The start state is q_0 , and we have to get to q_1 , since this is the only final state. Any number of a symbols can be consumed in q_0 , but a b symbol is required to transition to q_1 . Once there, any number of b symbols can be consumed, but an a symbol cannot. So the regular expression corresponding to the language defined by M_1 is a^*bb^* .

Computational properties of finite-state acceptors

The key computational question for finite-state acceptors is: how fast can we determine whether a string is accepted? For deterministic FSAs, this computation can be performed by Dijkstra's algorithm, with time complexity $\mathcal{O}(V \log V + E)$, where V is the number of vertices in the FSA, and E is the number of edges (Cormen et al., 2009). Non-deterministic FSAs (NFSAs) can include multiple transitions from a given symbol and state. Any NSFA can be converted into a deterministic FSA, but the resulting automaton may have a number of states that is exponential in the number of size of the original NFSFA (Mohri et al., 2002).

8.1.2 Morphology as a regular language

Many words have internal structure, such as prefixes and suffixes that shape their meaning. The study of word-internal structure is the domain of **morphology**, of which there are two main types:

- **Derivational morphology** describes the use of affixes to convert a word from one grammatical category to another (e.g., from the noun *grace* to the adjective *graceful*), or to change the meaning of the word (e.g., from *grace* to *disgrace*).
- **Inflectional morphology** describes the addition of details such as gender, number, person, and tense (e.g., the *-ed* suffix for past tense in English).

Morphology is a rich topic in linguistics, deserving of a course in its own right.¹ The focus here will be on the use of finite-state automata for morphological analysis. The

¹A good starting point would be a chapter from a linguistics textbook (e.g., Akmajian et al., 2010; Bender, 2013). A key simplification in this chapter is the focus on affixes as the sole method of derivation and inflection. English makes use of affixes, but also incorporates **apophony**, such as the inflection of *foot* to *feet*. In semitic languages such as Arabic and Hebrew, morphology involves a template-based system, in which roots

current section deals with derivational morphology; inflectional morphology is discussed in § 8.1.4.

Suppose that we want to write a program that accepts only those words that are constructed in accordance with the rules of English derivational morphology:

- (8.1) grace, graceful, gracefully, *gracelyful
- (8.2) disgrace, *ungrace, disgraceful, disgracefully
- (8.3) allure, *allureful, alluring, alluringly
- (8.4) fairness, unfair, *disfair, fairly

(Recall that the asterisk indicates that a linguistic example is judged unacceptable by fluent speakers of a language.) These examples cover only a tiny corner of English derivational morphology, but a number of things stand out. The suffix *-ful* converts the nouns *grace* and *disgrace* into adjectives, and the suffix *-ly* converts adjectives into adverbs. These suffixes must be applied in the correct order, as shown by the unacceptability of **gracelyful*. The *-ful* suffix only works for some words, as shown by the use of *alluring* as the adjectival form of *allure*. Other changes are made with prefixes, such as the derivation of *disgrace* from *grace*, which roughly corresponds to a negation; however, the negation of *fair* is performed with the *un-* prefix instead. Finally, while the first three examples suggest that the direction of derivation is noun \rightarrow adjective \rightarrow adverb, the example of *fair* suggests that the adjective can also be the base form, with the *-ness* suffix performing the conversion to a noun.

These examples lead to a computational question: can we build a computer program that accepts only well-formed English words, and rejects all others? This might at first seem trivial to solve with a brute-force attack: simply make a dictionary of all valid English words, as in a spell-checker. But such an approach fails to account for morphological **productivity** — the applicability of existing morphological rules to new words and names, such as *Trump* to *Trumpy* and *Trumpkin*, and *Clinton* to *Clintonian* and *Clintonite*. We need an approach that represents morphological rules explicitly, and for this we will try a finite state acceptor.

The dictionary approach can be implemented as a finite state acceptor, with the vocabulary Σ equal to the vocabulary of English, and a transition from the start state to the accepting state for each word. But this would of course fail to generalize beyond the original vocabulary, and would not capture anything about the **morphotactic** rules that govern derivations from new words. The first step towards a more general approach is shown in Figure 8.2, which is the state diagram for a finite-state transducer in which the vocabulary

are triples of consonants (e.g., *ktb*), and words are created by adding consonants: *kataba* (Arabic: he wrote), *kutub* (books), *maktab* (desk). For more detail on morphology, see texts from Haspelmath and Sims (2013) and Lieber (2015).

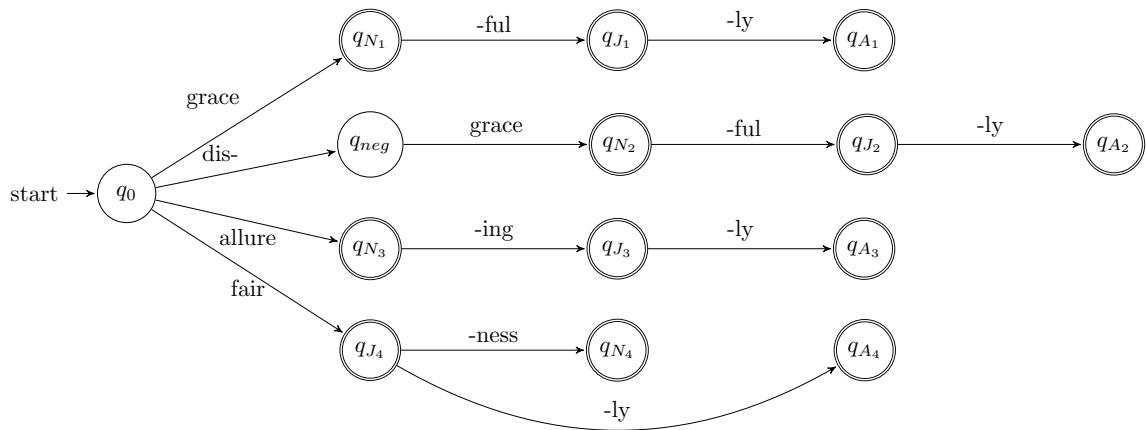


Figure 8.2: A finite-state acceptor for a fragment of English derivational morphology. Each path represents possible derivations from a single root form.

consists of **morphemes**, which include **stems** (e.g., *grace*, *allure*) and **affixes** (e.g., *dis-*, *-ing*, *-ly*). This finite-state acceptor consists of a set of paths leading away from the start state, with derivational affixes added along the path. Except for q_{neg} , the states on these paths are all final, so the FSA will accept *disgrace*, *disgraceful*, and *disgracefully*, but not *dis-*.

This FSA can be **minimized** to the form shown in Figure 8.3, which makes the generality of the finite-state approach more apparent. For example, the transition from q_0 to q_{J_2} can be made to accept not only *fair* but any single-morpheme (**monomorphemic**) adjective that takes *-ness* and *-ly* as suffixes. In this way, the finite-state acceptor can easily be extended: as new word stems are added to the vocabulary, their derived forms will be accepted automatically. Of course, this FSA would still need to be extended considerably to cover even this small fragment of English morphology. As shown by cases like *music* \rightarrow *musical*, *athlete* \rightarrow *athletic*, English includes several classes of nouns, each with its own rules for derivation.

The FSAs shown in Figure 8.2 and 8.3 accept *allureing*, not *alluring*. This reflects a distinction between morphology — the question of which morphemes and in what order — and **orthography** — the question of how the morphemes are rendered in written language. Just as orthography requires dropping the *e* preceding the *-ing* suffix, **phonology** imposes a related set of constraints on how words are rendered in speech. As we will see soon, these issues are handled through **finite-state transducers**, which are finite-state automata that take inputs and produce outputs.

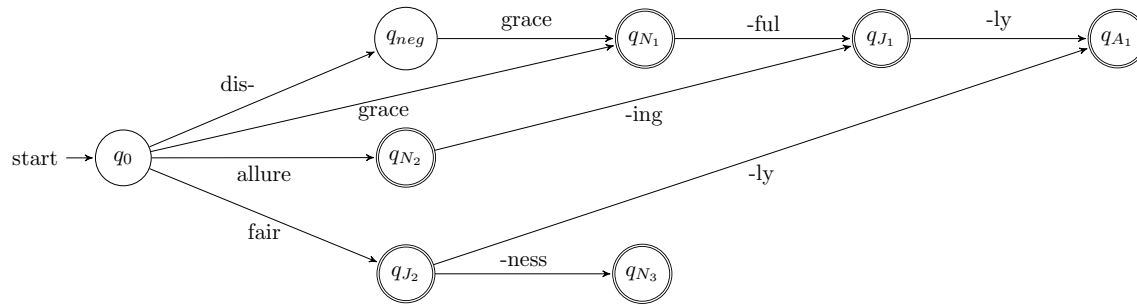


Figure 8.3: Minimization of the finite-state acceptor shown in Figure 8.2.

8.1.3 Weighted finite-state acceptors

According to the FSA treatment of morphology, every word is either in or out of the language, with no wiggle room. Perhaps you agree that *musicky* and *fishful* are not valid English words; but if forced to choose, you probably find *a fishful stew* or *a musicky tribute* preferable to *behaving disgracelyful*. Rather than asking whether a word is acceptable, we might like to ask how acceptable it is. Aronoff (1976, page 36) puts it another way: “Though many things are possible in morphology, some are more possible than others.” But finite state acceptors gives us no way to express preferences among technically valid choices.

Weighted finite-state acceptors (WFSAs) are generalizations of FSAs, in which each accepting path is assigned a score, which is computed from the transitions, the initial state, and the final state. Formally, a weighted finite-state acceptor $M = \langle Q, \Sigma, \lambda, \rho, \delta \rangle$ consists of:

- a finite set of states $Q = \{q_0, q_1, \dots, q_n\}$;
- a finite alphabet Σ of input symbols;
- an initial weight function, $\lambda : Q \rightarrow \mathbb{R}$;
- a final weight function $\rho : Q \rightarrow \mathbb{R}$;
- a transition function $\delta : Q \times \Sigma \times Q \rightarrow \mathbb{R}$.

Thus, WFSAs depart from the FSA formalism in three ways: every state can be an initial state, with score $\lambda(q)$; every state can be an accepting state, with score $\rho(q)$; transitions are possible between any pair of states on any input, with a score $\delta(q_i, \omega, q_j)$. Nonetheless, FSAs can be viewed as a special case: for any FSA M we can build an equivalent WFSa by setting $\lambda(q) = \infty$ for all $q \neq q_0$, $\rho(q) = \infty$ for all $q \notin F$, and $\delta(q_i, \omega, q_j) = \infty$ for all transitions $\{(q_i, \omega) \rightarrow q_j\}$ that are not permitted by the transition function of M .

The total score for any path $\pi = t_1, t_2, \dots, t_N$ is equal to the sum of these scores,

$$d(\pi) = \lambda(\text{from-state}(t_1)) + \sum_n^N \delta(t_n) + \rho(\text{to-state}(t_N)). \quad [8.5]$$

A **shortest-path algorithm** is used to find the minimum-cost path through a WFSAs for 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).²

N-gram language models as WFSAs

§ 5.1 introduced **n -gram language models**, in which the probability of a sequence of tokens w_1, w_2, \dots, w_M is modeled as,

$$p(w_1, \dots, w_M) \approx \prod_m^M p_n(w_m \mid w_{m-1}, \dots, w_{m-n+1}). \quad [8.6]$$

N -gram language models can be modeled as WFSAs. First consider a unigram language model. We need only a single state q_0 , with transition scores $\delta(q_0, \omega, q_0) = \log p_1(\omega)$. The initial and final scores can be set to zero. Then the path score for w_1, w_2, \dots, w_M is equal to,

$$0 + \sum_m^M \delta(q_0, w_m, q_0) + 0 = \sum_m^M \log p_1(w_m). \quad [8.7]$$

For an n -gram language model with $n > 1$, we need probabilities that condition on the past history. For example, in a bigram language model, the transition weights must represent $\log p_2(w_m \mid w_{m-1})$. This means that the transitions scoring function has to somehow “remember” the previous word or words. This can be done by adding more states: to model the bigram probability $p_2(w_m \mid w_{m-1})$, we need a state for every possible w_{m-1} — a total of V states. The construction indexes each state q_i by a context event $w_{m-1} = i$. The weights are then assigned as follows:

$$\begin{aligned} \delta(q_i, \omega, q_j) &= \begin{cases} \log \Pr_2(w_m = j \mid w_{m-1} = i), & \omega = j \\ \infty, & \omega \neq j \end{cases} \\ \lambda(q_i) &= \log \Pr_2(w_1 = i \mid w_0 = \diamond) \\ \rho(q_i) &= \log \Pr_2(w_{M+1} = \blacklozenge \mid w_M = i). \end{aligned}$$

²Shortest-path algorithms find the path with the minimum cost. In many cases, the path weights are log probabilities, so we want the path with the maximum score, which can be accomplished by making each local score into a **negative** log-probability. The remainder of this section will refer to **best-path algorithms**, which are assumed to “do the right thing.”

The transition function is designed to ensure that the context is recorded accurately: we can move to state j on input ω only if $\omega = j$; otherwise, transitioning to state j is forbidden by the weight of $-\infty$. The initial weight function $\lambda(q_i)$ is the log probability of receiving i as the first token, and the final weight function $\rho(q_i)$ is the log probability of receiving an “end-of-string” token after observing $w_M = i$. [todo: figure]

*Semiring Weighted Finite State Acceptors

The n -gram language model WFSAs is deterministic: each input has exactly one accepting path, for which the WFSAs computes a score. In non-deterministic WFSAs, a given input may have multiple accepting paths, and in some applications, we are interested in the score for the input across all such paths. Such aggregate scores can be computed by generalizing WFSAs with **semiring notation**, first introduced in § 6.4.4.

Let $d(\pi)$ represent the total score for path $\pi = t_1, t_2, \dots, t_N$, which is computed as,

$$d(\pi) = \lambda(\text{from-state}(t_1)) \otimes \delta(t_1) \otimes \delta(t_2) \otimes \dots \otimes \delta(t_N) \otimes \rho(\text{to-state}(t_N)). \quad [8.8]$$

This is a generalization of Equation 8.5 to semiring notation, using the semiring multiplication operator \otimes in place of addition.

Now let $s(\omega)$ represent the total score for all paths $\Pi(\omega)$ that consume input ω ,

$$s(\omega) = \bigoplus_{\pi \in \Pi(\omega)} d(\pi). \quad [8.9]$$

Here, semiring addition (\oplus) is used to combine the scores of multiple paths.

The generalization to semirings covers a number of useful special cases. In the log-probability semiring, multiplication is defined as $\log p(x) \otimes \log p(y) = \log p(x) + \log p(y)$, and addition is defined as $\log p(x) \oplus \log p(y) = \log(p(x) + p(y))$. Thus, $s(\omega)$ represents the log-probability of accepting input ω , marginalizing over all paths $\pi \in \Pi(\omega)$. In the **boolean semiring**, the \otimes operator is logical conjunction, and the \oplus operator is logical disjunction. This reduces to the special case of unweighted finite-state acceptors, where the score $s(\omega)$ is a boolean indicating whether there exists any accepting path for ω . In the **tropical semiring**, the \oplus operator is a maximum, so the resulting score is the score of the best-scoring path through the WFSAs. The `OpenFST` toolkit uses semirings and polymorphism to implement general algorithms for weighted finite-state automata (Allauzen et al., 2007).

*Interpolated n -gram language models

Recall from § 5.2.3 that an **interpolated n -gram language model** combines the probabilities from multiple n -gram models. For example, an interpolated bigram language model

(c) Jacob Eisenstein 2014-2017. Work in progress.

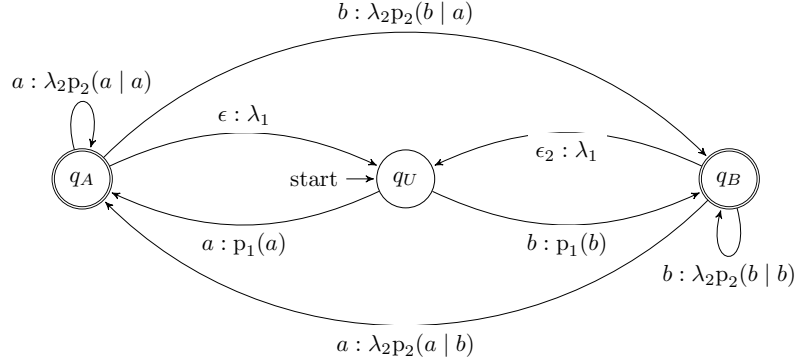


Figure 8.4: WFS A implementing an interpolated bigram/unigram language model, on the alphabet $\Sigma = \{a, b\}$. For simplicity, the WFS A is constrained to force the first token to be generated from the unigram model, and does not model the emission of the end-of-sequence token.

computes probability,

$$\hat{p}(w_m \mid w_{m-1}) = \lambda_1 p_1(w_m) + \lambda_2 p_2(w_m \mid w_{m-1}), \quad [8.10]$$

with \hat{p} indicating the interpolated probability, p_2 indicating the bigram probability, and p_1 indicating the unigram probability. We set $\lambda_2 = (1 - \lambda_1)$ so that the probabilities sum to one.

Knight and May (2009) describe how to implement an interpolated bigram language model using a non-deterministic WFS A. The basic idea is shown in Figure 8.4. In the WFS A implementation of the bigram language model, we have one state for each element in the vocabulary — in this case, the states q_A and q_B — which capture the contextual conditioning in the bigram probabilities. To model unigram probabilities, we add an additional state q_U , which “forgets” the context. Transitions out of q_U involve unigram probabilities, $p_1(a)$ and $p_2(b)$; transitions into q_U emit the empty symbol ϵ , and have probability λ_1 , reflecting the interpolation weight for the unigram model. The interpolation weight for the bigram model is included directly in the weight of the transition $q_A \rightarrow q_B$.

The epsilon transitions into q_U make this WFS A non-deterministic. Consider the score for the sequence (a, b, b) . The initial state is q_U , so the symbol a is generated with score $p_1(a)$ ³ Next, we can generate b from the unigram model by taking the transition $q_A \rightarrow q_B$, with score $\lambda_2 p_2(b \mid a)$. Alternatively, we can take a transition back to q_U with score λ_1 , and then emit b from the unigram model with score $p_1(b)$. To generate the final b token,

³We could model the sequence-initial bigram probability $p_2(a \mid \square)$, but for simplicity the WFS A does not admit this possibility, which would require another state.

we face the same choice: emit it directly from the self-transition to q_B , or transition to q_U first.

The total score for the sequence (a, b, b) is the semiring sum over all accepting paths,

$$\begin{aligned}
 s(a, b, b) = & (p_1(a) \otimes \lambda_2 p_2(b \mid a) \otimes \lambda_2 p(b \mid b)) \\
 & \oplus (p_1(a) \otimes \lambda_1 \otimes p_1(b) \otimes \lambda_2 p(b \mid b)) \\
 & \oplus (p_1(a) \otimes \lambda_2 p_2(b \mid a) \otimes p_1(b) \otimes p_1(b)) \\
 & \oplus (p_1(a) \otimes \lambda_1 \otimes p_1(b) \otimes p_1(b)) .
 \end{aligned} \tag{8.11}$$

Each line in Equation 8.11 represents the probability of a specific path through the WFSM. In the probability semiring, \otimes is multiplication, so that each path is the product of each transition weight, which are themselves probabilities. The \oplus operator is addition, so that the total score is the sum of the scores (probabilities) for each path. This corresponds to the probability under the interpolated bigram language model.

8.1.4 Finite-state transducers

Finite state acceptors can determine whether a string is in a regular language, and weighted finite state acceptors can compute a score for every string over a given alphabet. **Finite-state transducers** (FSTs) extend the formalism further, by adding an output symbol to each transition. Formally, we write $T = \langle Q, \Sigma, \Omega, \lambda, \rho, \delta \rangle$, with Ω representing an output vocabulary and the transition function $\delta : Q \times (\Sigma \cup \epsilon) \times (\Omega \cup \epsilon) \times Q \rightarrow \mathbb{R}$ mapping from states, input symbols, and output symbols to states. The remaining elements $(Q, \Sigma, \lambda, \rho)$ are identical to their definition in weighted finite state acceptors (§ 8.1.3). Thus, each path through the FST T transduces the input string into an output.

String edit distance

The **edit distance** between two strings s and t is a measure of how many basic operations are required to transform one string into another. There are several ways to compute edit distance, but one of the most popular is the **Levenshtein edit distance**, which counts the minimum number of insertions, deletions, and substitutions. This can be computed by a one-state weighted finite-state transducer, in which the input and output alphabets are identical. For simplicity, consider the alphabet $\Sigma = \Omega = [a, b]$. The edit distance can be computed by a one-state transducer with the following transitions,

$$\delta(q, a, a, q) = \delta(q, b, b, q) = 0 \tag{8.12}$$

$$\delta(q, a, b, q) = \delta(q, b, a, q) = 1 \tag{8.13}$$

$$\delta(q, a, \epsilon, q) = \delta(q, b, \epsilon, q) = 1 \tag{8.14}$$

$$\delta(q, \epsilon, a, q) = \delta(q, \epsilon, b, q) = 1. \tag{8.15}$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

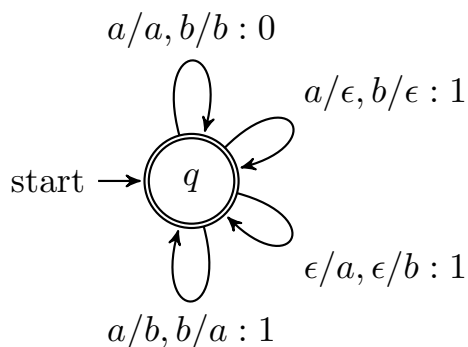


Figure 8.5: State diagram for the Levenshtein edit distance finite state transducer. The label $x/y : c$ indicates a cost of c for a transition with input x and output y .

The state diagram is shown in Figure 8.5.

For a given string pair, there are multiple paths through the transducer: the best-scoring path from *dessert* to *desert* involves a single deletion, for a total score of 1; the worst-scoring path has a score of 13. In edit distance, we are interested in the lowest score, but as we will see, in other applications, the score is computed by marginalizing over paths.

The Porter stemmer

The Porter (1980) stemming algorithm was discussed in chapter 3. It is a “lexicon-free” algorithm for stripping suffixes from English words, using a sequence of character-level rules. Each rule can be described by an unweighted finite-state transducer. The first rule is:

$-sses \rightarrow -ss$ e.g., *dresses* \rightarrow *dress* [8.16]

$-ies \rightarrow -i$ e.g., *parties* \rightarrow *parti* [8.17]

$-ss \rightarrow -ss$ e.g., *dress* \rightarrow *dress* [8.18]

$-s \rightarrow \epsilon$ e.g., *cats* \rightarrow *cat* [8.19]

The final two lines appear to conflict; they are meant to be interpreted as an instruction to remove a terminal $-s$ unless it is part of an $-ss$ ending. A state diagram to handle just these final two lines is shown in Figure 8.6. Be sure to understand how this finite-state transducer handles *cats*, *steps*, *bass*, and *basses*.

Inflectional morphology

In **inflectional morphology**, word **lemmas** are modified to add grammatical information such as tense, number, and case. For example, many English nouns are pluralized by the -

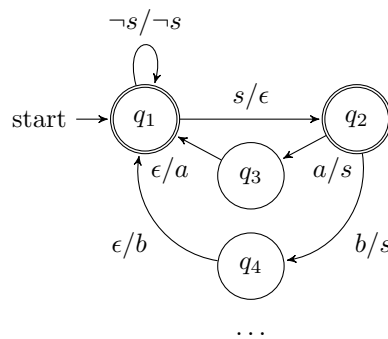


Figure 8.6: State diagram for final two lines of step 1a of the Porter stemming diagram. States q_3 and q_4 “remember” the observations a and b respectively; the ellipsis \dots represents additional states for each symbol in the input alphabet. The shorthand notation $\neg s/\neg s$ is not part of the FST formalism, but should be interpreted as a set of self-transition arcs for every input/output symbol except s .

infinitive	cantar (to sing)	comer (to eat)	vivir (to live)
yo (1st singular)	canto	como	vivo
tu (2nd singular)	cantas	comes	vives
él, ella, usted (3rd singular)	canta	come	vive
nosotros (1st plural)	cantamos	comemos	vivimos
vosotros (2nd plural, informal)	cantáis	coméis	vivís
ellos, ellas (3rd plural); ustedes (2nd plural)	cantan	comen	viven

Table 8.1: Spanish verb inflections for the present indicative tense. Each row represents a person and number, and each column is a regular example from a class of verbs, as indicated by the ending of the infinitive form.

s suffix, and many verbs are converted to past tense by the *-ed* suffix. English’s inflectional morphology is considerably simpler than many of the world’s languages. For example, Romance languages (derived from Latin) feature complex systems of verb suffixes which must agree with the person and number of the verb, as shown in Table 8.1.

The task of **morphological analysis** is to read a form like *canto*, and output an analysis like CANTAR+VERB+PRESIND+1P+SING, where +PRESIND describes the tense as present indicative, +1P indicates the first-person, and +SING indicates the singular number. The task of **morphological generation** is the reverse, going from CANTAR+VERB+PRESIND+1P+SING to *canto*. Finite-state transducers are an attractive solution, because they can solve both problems with a single model (Beesley and Karttunen, 2003). As an example, Figure 8.7

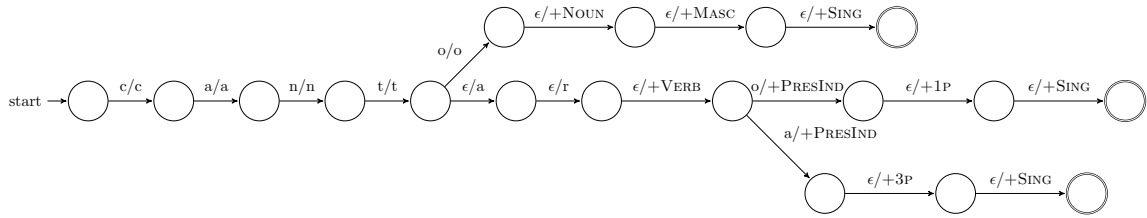


Figure 8.7: Fragment of a finite-state transducer for Spanish morphology. There are two accepting paths for the input *canto*: *canto*+NOUN+MASC+SING (masculine singular noun, meaning a song), and *cantar*+VERB+PRESIND+1P+SING (I sing). There is also an accepting path for *canta*, with output *cantar*+VERB+PRESIND+3P+SING (he/she sings).

shows a fragment of a finite-state transducer for Spanish inflectional morphology. The input vocabulary Σ corresponds to the set of letters used in Spanish spelling, and the output vocabulary Ω corresponds to these same letters, plus the vocabulary of morphological features (e.g., +SING, +VERB). In Figure 8.7, there are two paths that take *canto* as input, corresponding to the verb and noun meanings; the choice between these paths could be guided by a part-of-speech tagger. By **inversion**, the inputs and outputs for each transition are switched, resulting in a finite-state generator, capable of producing the correct **surface form** for any morphological analysis.

Finite-state morphological analyzers and other unweighted transducers can be designed by hand. The designer’s goal is to avoid **overgeneration** — accepting strings or making transductions that are not valid in the language — as well as **undergeneration** — failing to accept strings or transductions that are valid. For example, a pluralization transducer that does not accept *foot/feet* would undergenerate. Suppose we “fix” the transducer to accept this example, but as a side effect, it now accepts *boot/beet*; the transducer would then be said to overgenerate. A transducer that accepts *foot/foots* but not *foot/feet* would both overgenerate and undergenerate.

Finite-state composition

Designing finite-state transducers to capture the full range of morphological phenomena in any real language is a huge task. In computer science, **modularization** is a key technique for handling such issues — decompose a large and unwieldy problem into a set of subproblems, each of which will hopefully have a concise solution. In finite-state automata, this can be performed by **composition**: feeding the output of one transducer T_1 as the input to another transducer T_2 , written $T_2 \circ T_1$. Formally, if there exists some y such that $(x, y) \in T_1$ (meaning that T_1 produces output y on input x), and $(y, z) \in T_2$, then $(x, z) \in (T_2 \circ T_1)$. Because finite-state transducers are **closed** under composition, there is guaranteed to be a single finite-state transducer that $T_3 = T_2 \circ T_1$, which can be constructed as a machine with one state for each pair of states in T_1 and T_2 (Mohri et al.,

2002).

Example: Morphology and orthography In English morphology, the suffix *-ed* is added to signal the past tense for many verbs: *cook*→*cooked*, *want*→*wanted*, etc. However, English **orthography** dictates that this process cannot produce a spelling with consecutive *e*'s, so that *bake*→*baked*, not *bakeed*. A modular solution is to build separate transducers for morphology and orthography. The morphological transducer T_M transduces from *bake*+PAST to *bake+ed*, with the + symbol indicating a segment boundary. The input alphabet of T_M includes the lexicon of words and the set of morphological features; the output alphabet includes the characters *a-z* and the + boundary marker. Next, an orthographic transducer T_O is responsible for the transductions *cook+ed* → *cooked*, and *bake+ed* → *baked*. The input alphabet of T_O must be the same as the output alphabet for T_M , and the output alphabet is simply the characters *a-z*. The composition $(T_O \circ T_M)$ then transduces from *bake*+PAST to the spelling *baked*. The design of T_O is left as an exercise.

Example: Hidden Markov models Hidden Markov models (chapter 6) can be viewed as weighted finite-state transducers. Recall that a hidden Markov model defines a joint probability over words and tags, $p(\mathbf{w}, \mathbf{y})$, which can be computed as a path through a trellis structure. This trellis is itself a weighted finite-state **acceptor**, with edges between all adjacent nodes $q_{m-1,i} \rightarrow q_{m,j}$ on input $Y_m = j$. The edge weights are log-probabilities,

$$\delta(q_{m-1,i}, Y_m = j, q_{m,j}) = \log p(w_m, Y_m = j \mid Y_{m-1} = i) \quad [8.20]$$

$$= \log p(w_m \mid Y_m = j) + \log \Pr(Y_m = j \mid Y_{m-1} = i). \quad [8.21]$$

Because there is only one possible transition for each tag Y_m , this WFSA is deterministic. The score for any tag sequence $\{y_m\}_{m=1}^M$ is the sum of these log-probabilities, corresponding to the total log probability $\log p(\mathbf{w}, \mathbf{y})$.

Furthermore, the trellis can be constructed by the composition of simpler FSTs.

- First, take construct a “transition” transducer to represent a bigram probability model over tag sequences, T_T . This transducer is almost identical to the n -gram language model acceptor in § 8.1.3: there is one state for each tag, and the edge weights equal to the transition log-probabilities, $\delta(q_i, j, j, q_j) = \log \Pr(Y_m = j \mid Y_{m-1} = i)$. Note that T_T is a transducer, with identical input and output at each arc; this makes it possible to compose T_T with other transducers.
- Next, construct an “emission” transducer to represent the probability of words given tags, T_E . This transducer has only a single state, with arcs for each word/tag pair, $\delta(q_0, i, j, q_0) = \log \Pr(W_m = j \mid Y_m = i)$. The input vocabulary is the set of all tags, and the output vocabulary is the set of all words.

(c) Jacob Eisenstein 2014-2017. Work in progress.

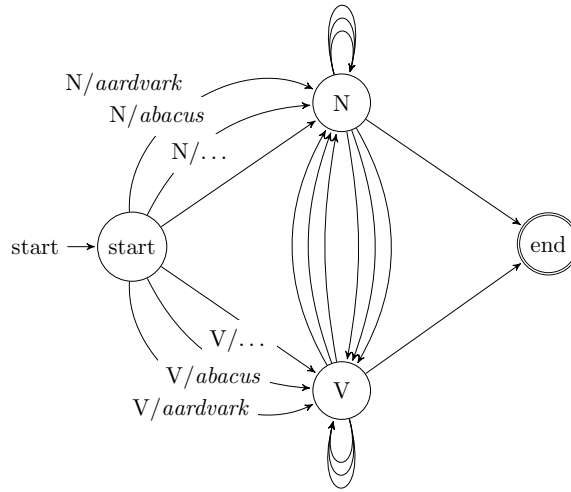


Figure 8.8: Finite-state transducer for hidden Markov models, with a small tagset of nouns and verbs. For each pair of tags (including self-loops), there is an edge for every word in the vocabulary. For simplicity, input and output are only shown for the edges from the start state. Weights are also omitted from the diagram; for each edge from q_i to q_j , the weight is equal to $\log p(w_m, Y_m = j \mid Y_{m-1} = i)$, except for edges to the end state, which are equal to $\log \Pr(Y_m = \blacklozenge \mid Y_{m-1} = i)$.

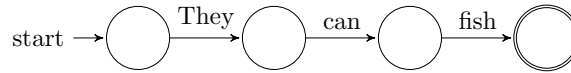


Figure 8.9: Chain finite-state acceptor for the input *They can fish*.

- The composition $T_E \circ T_T$ is a finite-state transducer with one state per tag, as shown in Figure 8.8. Each state has $V \times K$ outgoing edges, representing transitions to each of the K other states, with outputs for each of the V words in the vocabulary. The weights for these edges are equal to $\delta(q_i, Y_m = j, w_m, q_j) = \log p(w_m, Y_m = j \mid Y_{m-1} = i)$.
- Recall the **trellis** from § 6.3, a structure with $M \times K$ nodes, for each of the M words to be tagged and each of the K tags in the tagset. The trellis can be built by composition of $(T_E \circ T_T)$ against an unweighted **chain FSA** $M_A(w)$ that is specially constructed to accept only a given input w_1, w_2, \dots, w_M , shown in Figure 8.9. The trellis for input w is built from the composition $M_A(w) \circ (T_E \circ T_T)$. Composing with the unweighted $M_A(w)$ does not affect the edge weights from $(T_E \circ T_T)$, but it selects the subset of paths that generate the word sequence w .

(c) Jacob Eisenstein 2014-2017. Work in progress.

8.1.5 *Learning weighted finite-state automata

In generative models such as n -gram language models and hidden Markov models, the edge weights correspond to log probabilities that can be obtained from relative frequency estimation. However, in other cases, we wish to learn the edge weights from input/output pairs. This is difficult in non-deterministic finite-state automata, because we do not observe the specific arcs that are traversed in accepting the input, or in transducing from input to output. The path through the automaton is a **latent variable**.

chapter 4 presented one solution for learning with latent variables: expectation maximization (EM). This involves computing a distribution $q(\cdot)$ over the latent variable, and iterating between updates to this distribution and updates to the parameters — in this case, the arc weights. The **forward-backward algorithm** (§ 6.5.3) describes a dynamic program for computing a distribution over arcs in the trellis structure of a hidden Markov model, but this is a special case of the more general problem for finite-state automata. Eisner (2002) describes an **expectation semiring**, which enables the expected number of transitions across each arc to be computed through a semiring shortest-path algorithm. Alternative approaches for generative models include Markov Chain Monte Carlo (Chiang et al., 2010) and spectral learning (Balle et al., 2011).

Further afield, we can take a perceptron-style approach, with each arc corresponding to a feature. The classic perceptron update would update the weights by subtracting the difference between the feature vector corresponding to the predicted path and the feature vector corresponding to the correct path. Since the path is not observed, we resort to a **hidden variable perceptron**. The model is described formally in § 11.4, but the basic idea is to compute an update from the difference between the features from the predicted path and the features for the best-scoring path that generates the correct output.

8.2 Context-free languages

Beyond the restricted class of regular languages lie the context-free languages. An example of a language that is context-free but not finite-state is the language of arithmetic expressions with balanced parentheses. Intuitively, to accept only strings in this language, an FSA would have to “count” the number of left parentheses, and make sure that they are balanced against the number of right parentheses. An arithmetic expression can be arbitrarily long, yet by definition an FSA has a finite number of states. Thus, for any FSA, there will be a string that with too many parentheses to count. More formally, the **pumping lemma** is a proof technique for showing that languages are not regular. It is typically demonstrated for the simpler case $a^n b^n$, the language of strings containing a sequence of a ’s, and then an equal-length sequence of b ’s.⁴

⁴Details of the proof can be found in introductory computer science theory textbooks (e.g., Sipser, 2012).

		the dog			
	the cat	the dog	chased		
the goat	the cat	the dog	chased	kissed	
		...			

Figure 8.10: Three levels of center embedding

There are at least two arguments for the relevance of non-regular formal languages to linguistics. First, there are natural language phenomena that are argued to be isomorphic to $a^n b^n$. For English, the classic example is **center embedding**, shown in Figure 8.10. The initial expression *the dog* specifies a single dog. Embedding this expression into *the cat ... chased* specifies a particular cat — the one chased by the dog. This cat can then be embedded again to specify a goat, in the less felicitous but arguably grammatical expression, *the goat the cat the dog chased kissed*, which refers to the goat who was kissed by the cat which was chased by the dog. Chomsky (1957) argues that to be grammatical, a center-embedded construction must be balanced: if it contains n noun phrases (e.g., *the cat*), they must be followed by exactly $n - 1$ verbs. An FSA that could recognize such expressions would also be capable of recognizing the language $a^n b^n$. Because we can prove that no FSA exists for $a^n b^n$, no FSA can exist for center embedded constructions either. English includes center embedding, and so the argument goes, English grammar as a whole cannot be regular.⁵

A more practical argument for moving beyond regular languages is **modularity**. Many linguistic phenomena — especially in syntax — involve constraints that apply at long distance. Consider the problem of determiner-noun number agreement in English: we can say *the coffee* and *these coffees*, but not **these coffee*. By itself, this is easy enough to model in an FSA. However, fairly complex modifying expressions can be inserted between the determiner and the noun:

(8.5) the burnt coffee

(8.6) the badly-ground coffee

(8.7) the burnt and badly-ground Italian coffee

(8.8) these burnt and badly-ground Italian coffees

(8.9) *these burnt and badly-ground Italian coffee

⁵The claim that arbitrarily deep center-embedded expressions are grammatical has drawn skepticism. Corpus evidence shows that embeddings of depth greater than two are exceedingly rare (Karlsson, 2007), and that embeddings of depth greater than three are completely unattested. If center-embedding is capped at some finite depth, then it is regular.

Again, an FSA can be designed to accept modifying expressions such as *burnt and badly-ground Italian*. Let's call this FSA F_M . To reject the final example, a finite-state acceptor must somehow “remember” that the determiner was plural when it reaches the noun *coffee* at the end of the expression. The only way to do this is to make two identical copies of F_M : one for singular determiners, and one for plurals. While this is possible in the finite-state framework, it is inconvenient — especially in languages where more than one attribute of the noun is marked by the determiner. **Context-free languages** are a superset of regular languages, which include cases such as $a^n b^n$, and which enable modularity even in the case of long-range dependencies such as determiner-noun agreement.

8.2.1 Context-free grammars

Context-free languages are specified by **context-free grammars (CFGs)**, which are tuples $\langle N, \Sigma, R, S \rangle$ consisting of:

- a finite set of **non-terminals** N ;
- a finite alphabet Σ of **terminal symbols**;
- a set of **production rules** R , each of the form $A \rightarrow \beta$, where $A \in N$ and $\beta \in (\Sigma \cup N)^*$;
- a designated start symbol S .

In the production rule $A \rightarrow \beta$, the left-hand side (LHS) A must be a non-terminal; the right-hand side (RHS) can be a sequence of terminals or non-terminals, $\{n, \sigma\}^*$, $n \in N, \sigma \in \Sigma$. A given non-terminal can appear on the left-hand side of many production rules. The non-terminal A on the left-hand side can also appear on the right-hand side; such recursive productions are analogous to self-loops in finite-state automata. Context-free grammars are “context-free” in the sense that the production rule depends only on the LHS, and not on its ancestors or neighbors; this is analogous to Markov property of finite-state automata, in which the behavior at each step depends only on the current state, on not on the path by which that state was reached.

A **derivation** τ is a sequence of steps from the start symbol S to a surface string $w \in \Sigma^*$, which is the **yield** of the derivation. A string w is in a context-free language if there is some derivation from S yielding w . **Parsing** is the problem of finding a derivation for a string in a grammar. Algorithms for parsing are described in chapter 9.

Context-free grammars are analogous to regular expressions, in that they define the language but not the computation necessary to recognize it. The context-free analogues to finite-state acceptors are **pushdown automata**, a theoretical model of computation in which input symbols can be pushed onto a stack with potentially infinite depth (Sipser, 2012).

$$\begin{aligned}
S &\rightarrow S \text{ OP } S \mid \text{NUM} \\
\text{OP} &\rightarrow + \mid - \mid \times \mid \div \\
\text{NUM} &\rightarrow \text{NUM DIGIT} \mid \text{DIGIT} \\
\text{DIGIT} &\rightarrow 0 \mid 1 \mid 2 \mid \dots \mid 9
\end{aligned}$$

Figure 8.11: A context free grammar for arithmetic expressions

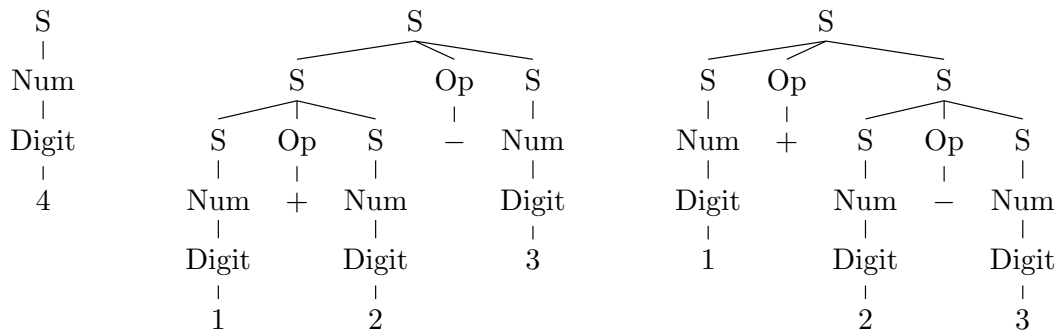


Figure 8.12: Some example derivations from the arithmetic grammar in Figure 8.11

Example

Figure 8.11 shows a context-free grammar for arithmetic expressions such as $(1 + 2) \div (3 \times 4)$. In this grammar, the terminal symbols include the digits $\{1, 2, \dots, 9\}$ and the operators $\{+, -, \times, \div\}$. The rules include the $|$ symbol, a notational convenience that makes it possible to specify multiple right-hand sides on a single line: the statement $A \rightarrow x \mid y$ defines *two* productions, $A \rightarrow x$ and $A \rightarrow y$. This grammar is recursive: the non-terminals S and NUM can produce themselves.

Derivations are typically shown as trees, with production rules applied from the top to the bottom. The tree on the left in Figure 8.12 describes the derivation of a single digit, through the sequence of productions $S \rightarrow \text{EXPR} \rightarrow \text{NUM} \rightarrow \text{DIGIT} \rightarrow 4$ (these are all **unary productions**, because the right-hand side contains a single element). The other two trees in Figure 8.12 show alternative derivations of the string $1 + 2 - 3$. The existence of multiple derivations for a string indicates that the grammar is **ambiguous**.

Context-free derivations can also be written out according to the pre-order tree traversal.

sal. For the two derivations of $1 + 2 - 3$ in Figure 8.12, the notation is:

$$(S (S (S (Num (Digit 1))) (Op +) (S (Num (Digit 2)))) (Op -) (S (Num (Digit 3)))) \quad [8.22]$$

$$(S (S (Num (Digit 1))) (Op +) (S (Num (Digit 2)) (Op -) (S (Num (Digit 3))))). \quad [8.23]$$

Grammar equivalence and Chomsky Normal Form

A single context-free language may be expressed by more than one context-free grammar: as a trivial example, the language $a^n b^n$ ($n > 0$) can be recognized by either of the following grammars:

$$\begin{aligned} S &\rightarrow aSb \mid ab \\ S &\rightarrow aSb \mid aabb \mid ab \end{aligned}$$

Two grammars are **weakly equivalent** if they generate the same strings. Two grammars are **strongly equivalent** if they generate the same strings via the same derivations. (The grammars above are only weakly equivalent.)

In **Chomsky Normal Form (CNF)**, the right-hand side of every production includes either two non-terminals, or a single terminal symbol:

$$\begin{aligned} A &\rightarrow BC \\ A &\rightarrow a \end{aligned}$$

All CFGs can be converted into a CNF grammar that is weakly equivalent. To convert a grammar into CNF, we first address productions that have more than two non-terminals on the RHS by creating new “dummy” non-terminals. For example, if we have the production,

$$W \rightarrow X Y Z, \quad [8.24]$$

we can replace it with two productions,

$$W \rightarrow X W \backslash X \quad [8.25]$$

$$W \backslash X \rightarrow Y Z. \quad [8.26]$$

In these productions, $W \backslash X$ is a new dummy non-terminal, indicating a phrase that would be W if its left neighbor is an X . This transformation **binarizes** the grammar, which is critical for efficient bottom-up parsing (chapter 9). Productions whose right-hand side contains a mix of terminal and non-terminal symbols can be replaced in a similar fashion.

Unary non-terminal productions $A \rightarrow B$ are replaced as follows: identify all productions $B \rightarrow \alpha$, and add $A \rightarrow \alpha$ to the grammar. For example, in the grammar described in Figure 8.11, we would replace $NUM \rightarrow DIGIT$ with $NUM \rightarrow 1 \mid 2 \mid \dots \mid 9$. However, we keep the production $NUM \rightarrow NUM DIGIT$, which is a valid binary production.

(c) Jacob Eisenstein 2014-2017. Work in progress.

8.2.2 Natural language syntax as a context-free language

Context-free grammars are widely used to represent **syntax**, which is the system of constraints and preferences that determine whether an utterance is judged to be grammatical. If this representation were perfectly faithful, then a natural language such as English could be transformed into a formal language, consisting of exactly the (infinite) set of strings that would be judged to be grammatical by a fluent English speaker. We could then build parsing software that would automatically determine if a given utterance were grammatical.⁶

Contemporary theories generally do *not* consider natural languages to be context-free, yet context-free grammars are widely used in automated natural language parsing. The reason is that context-free representations strike a good balance: it is possible to represent a broad range of syntactic phenomena in context-free form, and it is possible to parse context-free languages efficiently. This section therefore describes how to handle a core fragment of English syntax in context-free form, following the conventions of the **Penn Treebank** (PTB; Marcus et al., 1993), a large-scale annotation of English language syntax. The generalization to “mildly” context sensitive languages is discussed in § 8.3.

The Penn Treebank annotation is a **phrase-structure grammar** of English. This means that sentences are broken down into **constituents**, which are contiguous sequences of words that function as coherent units for the purpose of linguistic analysis. Constituents can be:

- moved around the sentence

(8.10) Abigail gave (her brother) (a fish).

(8.11) Abigail gave (a fish) to (her brother).

- substituted

(8.12) Max thanked (his older sister).

(8.13) Max thanked (her).

- coordinated

(8.14) (Abigail) and (her younger brother) bought a fish.

(8.15) Abigail (bought a fish) and (gave it to Max).

(8.16) Abigail (bought) and (greedily ate) a fish.

⁶As with morphology, it is impossible to do justice to natural language syntax here; the reader is encouraged to consult a text on linguistics (Akmajian et al., 2010; Bender, 2013).

These examples argue for units such as *her brother* and *bought a fish* to be treated as constituents. Other sequences of words in these examples, such as *Abigail gave* and *brother a fish*, cannot be moved, substituted, and coordinated in these ways. In phrase-structure grammar, constituents are nested, so that *the senator from New Jersey* contains the constituent *from New Jersey*, which in turn contains *New Jersey*. The sentence itself is the maximal constituent; each word is a minimal constituent, derived from a unary production from a part-of-speech tag. Between part-of-speech tags and sentences are **phrases**. In phrase-structure grammar, phrases are typically defined as the **projection** of a **head word**: for example, a **noun phrase** corresponds to a noun and the group of words that modify it, such as *her younger brother*; a **verb phrase** includes the verb and its modifiers, such as *bought a fish* and *greedily ate it*.

In context-free grammars, constituents correspond to the non-terminals. Grammar design involves choosing the right set of non-terminals. Fine-grained non-terminals make it possible to represent more fine-grained linguistic phenomena. For example, by distinguishing singular and plural noun phrases, it is possible to design a grammar rule that enforces subject-verb agreement in English. However, enforcing subject-verb agreement is considerably more complicated in languages like Spanish, where the verb must agree in both person and number with subject. In general, grammar designers must trade off between **overgeneration** — a grammar that permits ungrammatical sentences — and **undergeneration** — a grammar that fails to permit grammatical sentences. Furthermore, if the grammar is to support manual annotation of syntactic structure, it must be simple enough to annotate efficiently.

8.2.3 A phrase-structure grammar for English

To better understand how phrase-structure grammar works, let's consider the specific case of the Penn Treebank grammar of English. The main phrase categories in the Penn Treebank (PTB) are based on the main part-of-speech classes: noun phrase (NP), verb phrase (VP), prepositional phrase (PP), adjectival phrase (ADJP), and adverbial phrase (ADVP). The top-level category is S, which conveniently stands in for both “sentence” and the “start” symbol. **Complement clauses** (e.g., *She said that she liked spicy food*) are represented by the non-terminal SBAR. The terminal symbols in the grammar are individual words, which come from unary productions from part-of-speech tags (the PTB tagset is described in § 7.1).

This section explores the productions from the major phrase-level categories, explaining how to generate individual tag sequences. The production rules are approached in a “theory-driven” manner: first the syntactic properties of each phrase type are described, and then some of the necessary production rules are listed. But it is important to keep in mind that the Penn Treebank was produced in a “data-driven” manner. After the set of non-terminals was specified, annotators were free to analyze each sentence in what-

ever way seemed most linguistically accurate, subject to some high-level guidelines. The grammar of the Penn Treebank is simply the set of productions that were required to analyze the several million words of the corpus. By design, the grammar does not exclude ungrammatical sentences.

Sentences

The most common production rule for sentences is,

$$S \rightarrow NP VP \quad [8.27]$$

which accounts for simple sentences like *Abigail ate the kimchi* — as we will see, the direct object *the kimchi* is part of the verb phrase. But there are more complex forms of sentences as well:

$$S \rightarrow ADVP NP VP \quad \text{Unfortunately Abigail ate the kimchi.} \quad [8.28]$$

$$S \rightarrow S CC S \quad \text{Abigail ate the kimchi and Max had a burger.} \quad [8.29]$$

$$S \rightarrow VP \quad \text{Eat the kimchi.} \quad [8.30]$$

where ADVP is an adverbial phrase (e.g., *unfortunately*, *very unfortunately*) and CC is a coordinating conjunction (e.g., *and*, *but*).⁷

Noun phrases

Noun phrases refer to entities, real or imaginary, physical or abstract: *Asha*, *the steamed dumpling*, *parts and labor*, *nobody*, and *the rise of revolutionary syndicalism in the early twentieth century*. Noun phrase productions include “bare” nouns, which may optionally follow determiners, as well as pronouns:

$$NP \rightarrow NN \mid NNS \mid NNP \mid PRP \quad [8.31]$$

$$NP \rightarrow DET NN \mid DET NNS \mid DET NNP \quad [8.32]$$

The part-of-speech tags NN, NNS, and NNP refer to singular, plural, and proper nouns; PRP refers to personal pronouns, and DET refers to determiners. The grammar also contains terminal productions from each of these tags, e.g., $PRP \rightarrow I \mid you \mid we \mid \dots$

Noun phrases may be modified by adjectival phrases (ADJP; e.g., *the small Russian dog*) and numbers (CD; e.g., *the five pastries*), each of which may optionally follow a determiner:

$$NP \rightarrow ADJP NN \mid ADJP NNS \mid DET ADJP NN \mid DET ADJP NNS \quad [8.33]$$

$$NP \rightarrow CD NNS \mid DET CD NNS \mid \dots \quad [8.34]$$

⁷Notice that the grammar does not include the recursive production $S \rightarrow ADVP S$. It may be helpful to think about why this production would cause the grammar to overgenerate.

Some noun phrases include multiple nouns, such as *the liberation movement* and *an antelope horn*, necessitating additional productions:

$$\text{NP} \rightarrow \text{NN NN} \mid \text{NN NNS} \mid \text{DET NN NN} \mid \dots \quad [8.35]$$

These multiple noun constructions can be combined with adjectival phrases and cardinal numbers, leading to a large number of additional productions.

Recursive noun phrase productions include coordination, prepositional phrase attachment, subordinate clauses, and verb phrase adjuncts:

$$\text{NP} \rightarrow \text{NP Cc NP} \quad \text{e.g., the red and the black} \quad [8.36]$$

$$\text{NP} \rightarrow \text{NP PP} \quad \text{e.g., the President of the Georgia Institute of Technology} \quad [8.37]$$

$$\text{NP} \rightarrow \text{NP SBAR} \quad \text{e.g., the bicycle that I found outside} \quad [8.38]$$

$$\text{NP} \rightarrow \text{NP VP} \quad \text{e.g., a bicycle made of titanium} \quad [8.39]$$

These recursive productions are a major source of ambiguity, because the VP and PP non-terminals can also generate NP children. Thus, the *the President of the Georgia Institute of Technology* can be derived in two ways, as can *a bicycle made of titanium found outside*.

But aside from these few recursive productions, the noun phrase fragment of the Penn Treebank grammar is relatively flat, containing a large number of productions that go from NP directly to a sequence of parts-of-speech. If noun phrases had more internal structure, the grammar would need fewer rules; as we will see, this would make parsing faster and machine learning easier. Vadas and Curran (2011) propose to add additional structure in the form of a new non-terminal called a **nominal modifier** (NML), e.g.,

$$(8.17) \quad (\text{NP} (\text{NN crude}) (\text{NN oil}) (\text{NNS prices})) \quad (\text{PTB analysis})$$

$$(8.18) \quad (\text{NP} (\text{NML} (\text{NN crude}) (\text{NN oil})) (\text{NNS prices})) \quad (\text{NML-style analysis})$$

Another proposal is to treat the determiner as the head of a **determiner phrase** (DP; Abney, 1987). There are linguistic arguments for and against determiner phrases (Van Eynde, 2006, e.g.). From the perspective of context-free grammar, DPs enable more structured analyses of some constituents, e.g.,

$$(8.19) \quad (\text{NP} (\text{DT the}) (\text{JJ bald}) (\text{NN man})) \quad (\text{PTB analysis})$$

$$(8.20) \quad (\text{DP} (\text{DT the}) (\text{NP} (\text{JJ bald}) (\text{NN man}))) \quad (\text{DP-style analysis})$$

Verb phrases

Verb phrases describe actions, events, and states of being. The PTB tagset distinguishes several classes of verb inflections: base form (VB; *she likes to snack*), present-tense third-person singular (VBZ; *she snacks*), present tense but not third-person singular (VBP; *they*

snack), past tense (VBD; *they snacked*), present participle (VBG; *they are snacking*), and past participle (VBN; *they had snacked*).⁸ Each of these forms can constitute a verb phrase on its own:

$$VP \rightarrow VB \mid VBZ \mid VBD \mid VBN \mid VBG \mid VBP \quad [8.40]$$

More complex verb phrases can be formed by a number of recursive productions, including the use of coordination, modal verbs (MD; *she should snack*), and the infinitival *to* (TO):

$$VP \rightarrow MD \ VP \quad \text{She will snack} \quad [8.41]$$

$$VP \rightarrow VBD \ VP \quad \text{She had snacked} \quad [8.42]$$

$$VP \rightarrow VBZ \ VP \quad \text{She has been snacking} \quad [8.43]$$

$$VP \rightarrow VBN \ VP \quad \text{She has been snacking} \quad [8.44]$$

$$VP \rightarrow TO \ VP \quad \text{She wants to snack} \quad [8.45]$$

$$VP \rightarrow VP \ VP \quad \text{She buys and eats many snacks} \quad [8.46]$$

Each of these productions uses recursion, with VP appearing on the right-hand side. This enables the creation of very complex verb phrases, such as *She will have wanted to have been snacking*.

Transitive verbs take noun phrases as direct objects, and ditransitive verbs take two direct objects:

$$VP \rightarrow VBZ \ NP \quad \text{She teaches algebra} \quad [8.47]$$

$$VP \rightarrow VBG \ NP \quad \text{She has been teaching algebra} \quad [8.48]$$

$$VP \rightarrow VBD \ NP \ NP \quad \text{She taught her brother algebra} \quad [8.49]$$

These productions are *not* recursive, so a unique production is required for each verb part-of-speech. They also do not distinguish transitive from intransitive verbs, so the resulting grammar overgenerates examples like **She sleeps sushi* and **She learns Boyang algebra*. Sentences can also be direct objects:

$$VP \rightarrow VBZ \ S \quad \text{Asha wants to eat the kimchi} \quad [8.50]$$

$$VP \rightarrow VBZ \ SBAR \quad \text{Asha knows that Boyang eats the kimchi} \quad [8.51]$$

The first production overgenerates, licensing sentences like **Asha sees Boyang eats the kimchi*. This problem could be addressed by designing a more specific set of sentence non-terminals, indicating whether the main verb can be conjugated.

⁸It bears emphasis the principles governing this tagset design are entirely English-specific: VBP is a meaningful category only because English morphology distinguishes third-person singular from all person-number combinations.

Verbs can also be modified by prepositional phrases and adverbial phrases:

VP \rightarrow VBZ PP *She studies at night* [8.52]

VP \rightarrow VBZ ADVP *She studies intensively* [8.53]

VP \rightarrow ADVP VBG *She is not studying* [8.54]

Again, because these productions are not recursive, the grammar must include productions for every verb part-of-speech.

A special set of verbs, known as **copula**, can take **predicative adjectives** as direct objects:

VP \rightarrow VBZ ADJP *She is hungry* [8.55]

VP \rightarrow VBP ADJP *Success seems increasingly unlikely* [8.56]

The PTB does not have a special non-terminal for copular verbs, so this production generates non-grammatical examples such as **She eats tall*.

Particles (PRT as a phrase; RP as a part-of-speech) work to create phrasal verbs:

VP \rightarrow VB PRT *She told them to fuck off* [8.57]

VP \rightarrow VBD PRT NP *They gave up their ill-gotten gains* [8.58]

As the second production shows, particle productions are required for all configurations of verb parts-of-speech and direct objects.

Other constituents

The remaining constituents require far fewer productions. **Prepositional phrases** almost always consist of a preposition and a noun phrase,

PP \rightarrow IN NP *United States of America* [8.59]

PP \rightarrow TO NP *He gave his kimchi to Abigail* [8.60]

Similarly, complement clauses consist of a complementizer (usually a preposition, possibly null) and a sentence,

SBAR \rightarrow IN S *She said that it was spicy* [8.61]

SBAR \rightarrow S *She said it was spicy* [8.62]

Adverbial phrases are usually bare adverbs (ADVP \rightarrow RB), with a few exceptions:

ADVP \rightarrow RB RBR *They went considerably further* [8.63]

ADVP \rightarrow ADVP PP *They went considerably further than before* [8.64]

(c) Jacob Eisenstein 2014-2017. Work in progress.

The tag RBR is a comparative adverb.

Adjectival phrases extend beyond bare adjectives ($\text{ADJP} \rightarrow \text{JJ}$) in a number of ways:

$\text{ADJP} \rightarrow \text{RB JJ}$	<i>very hungry</i>	[8.65]
$\text{ADJP} \rightarrow \text{RBR JJ}$	<i>more hungry</i>	[8.66]
$\text{ADJP} \rightarrow \text{JJS JJ}$	<i>best possible</i>	[8.67]
$\text{ADJP} \rightarrow \text{RB JJR}$	<i>even bigger</i>	[8.68]
$\text{ADJP} \rightarrow \text{JJ CC JJ}$	<i>high and mighty</i>	[8.69]
$\text{ADJP} \rightarrow \text{JJ JJ}$	<i>West German</i>	[8.70]
$\text{ADJP} \rightarrow \text{RB VBN}$	<i>previously reported</i>	[8.71]

The tags JJR and JJS refer to comparative and superlative adjectives respectively.

All of these phrase types can be coordinated:

$\text{PP} \rightarrow \text{PP CC PP}$	<i>on time and under budget</i>	[8.72]
$\text{ADVP} \rightarrow \text{ADVP CC ADVP}$	<i>now and two years ago</i>	[8.73]
$\text{ADJP} \rightarrow \text{ADJP CC ADJP}$	<i>quaint and rather deceptive</i>	[8.74]
$\text{SBAR} \rightarrow \text{SBAR CC SBAR}$	<i>whether they want control</i> <i>or whether they want exports</i>	[8.75]

8.2.4 Grammatical ambiguity and weighted context-free grammars

Context-free parsing is useful not only because it determines whether a sentence is grammatical, but mainly because it breaks the sentence up into constituents, which are useful in applications such as information extraction (chapter 16) and sentence compression (Jing, 2000; Clarke and Lapata, 2008). However, the **ambiguity** of wide-coverage natural language grammars poses a serious problem for such potential applications. As an example, Figure 8.13 shows two possible analyses for the simple sentence *she eats sushi with chopsticks*, depending on whether the *chopsticks* are modifying *eats* or *sushi*. Realistic grammars license thousands or even millions of parses for individual sentences.

Weighted context-free grammars solve this problem by attaching weights to each production. The score of a derivation $\tau = \{(X \rightarrow \alpha)\}$ is then the sum of the weights of the productions,

$$s(\tau) = \sum_{(X \rightarrow \alpha) \in \tau} \psi(X \rightarrow \alpha), \quad [8.76]$$

where $\psi(X \rightarrow \alpha)$ is the score of the production $X \rightarrow \alpha$. The parsing problem is then,

$$\hat{\tau} = \operatorname{argmax}_{\tau: \text{yield}(\tau)=w} s(\tau). \quad [8.77]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

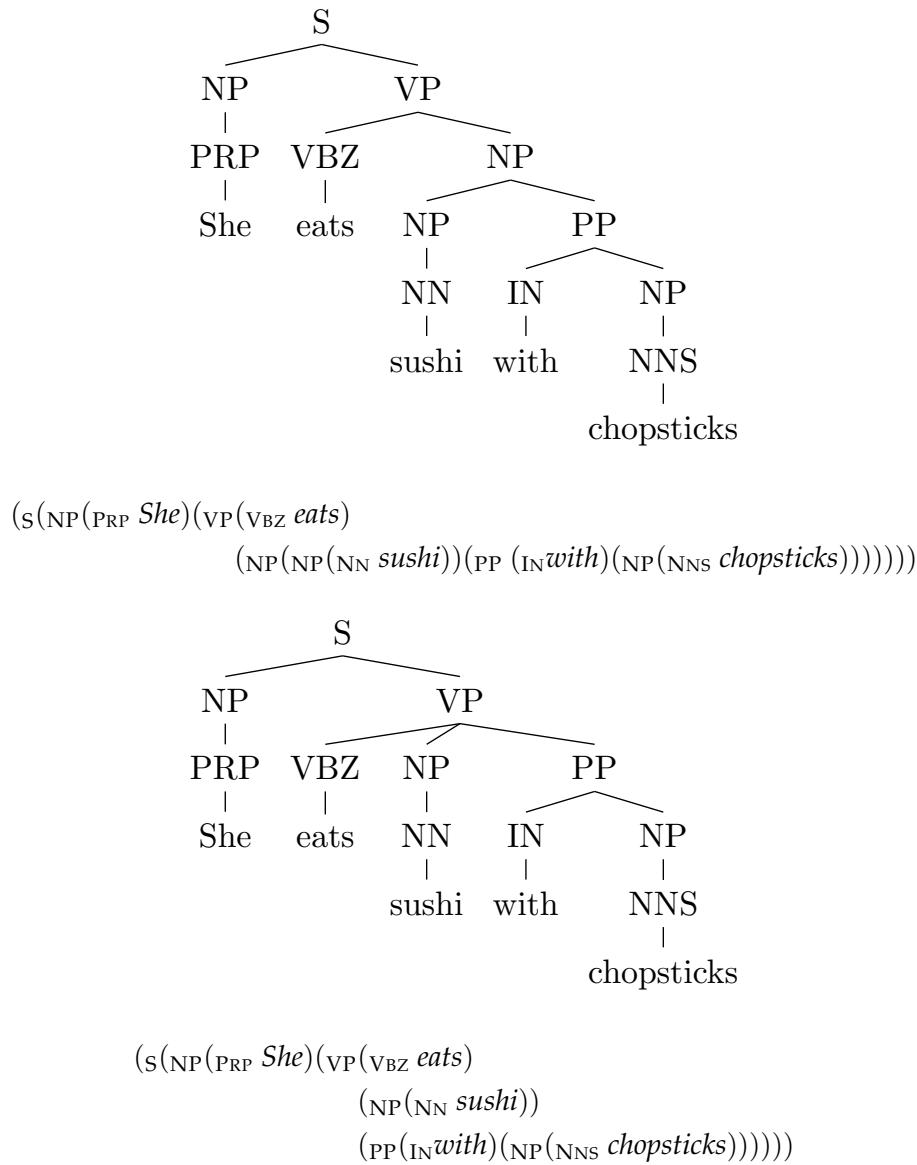


Figure 8.13: Two derivations of the same sentence, shown as both parse trees and bracketings

Parsing algorithms for weighted CFGs are described in chapter 9; in general, the time complexity is cubic in the length of the input. Given a labeled dataset, the weights can be set equal to log conditional probabilities, $\psi(X \rightarrow \alpha) = \log p(\alpha \mid X)$. Setting the weights in this way defines a **probabilistic context-free grammar** (PCFG), in which the total score of a derivation is equal to the log of the joint probability, $\log p(w, \tau)$. Alternatively, the weights can be learned using discriminative algorithms such as perceptron. For more details, see § 9.3.

8.3 *Mildly context-sensitive languages

Beyond context-free languages lie **context-sensitive languages**, in which the expansion of a non-terminal depends on its neighbors. In the general class of context-sensitive languages, computation becomes much more challenging: the membership problem for context-sensitive languages is PSPACE-complete. Since PSPACE contains the complexity class NP (problems that can be solved in polynomial time on a non-deterministic Turing machine), PSPACE-complete problems cannot be solved efficiently if $P \neq NP$. Thus, designing an efficient parsing algorithm for the full class of context-sensitive languages is probably hopeless.⁹

However, Joshi (1985) identifies a set of properties that define **mildly context-sensitive languages**, which are a strict subset of context-sensitive languages. Like context-free languages, mildly context-sensitive languages are efficiently parseable. However, the mildly context-sensitive languages include non-context-free languages, such as the “copy language” $\{ww \mid w \in \Sigma^*\}$ and the language $a^m b^n c^m d^n$. Both are characterized by **cross-serial dependencies**, linking symbols at long distance across the string.¹⁰ For example, in the language $a^n b^m c^n d^m$, each a symbol is linked to exactly one c symbol, regardless of the number of intervening b symbols.

Shieber (1985) describes a fragment of Swiss-German, in which sentences such as *we let the children help Hans paint the house* are realized by listing all nouns before all verbs, i.e., *we the children Hans the house let help paint*. Furthermore, each noun’s determiner is dictated by the noun’s **case marking** (the role it plays with respect to the verb). Using an argument that is analogous to the earlier discussion of center-embedding (§ 8.2), Shieber argues that these case marking constraints are a cross-serial dependency, homomorphic to $a^m b^n c^m d^n$, and therefore not context-free.

⁹If $PSPACE \neq NP$, then it contains problems that cannot be solved in polynomial time on a non-deterministic Turing machine; equivalently, solutions to these problems cannot even be checked in polynomial time (Arora and Barak, 2009).

¹⁰A further condition of the set of mildly-context sensitive languages is **constant growth**: if the strings in the language are arranged by length, the gap in length between any pair of adjacent strings is bounded by some language specific constant. This condition excludes languages such as $\{a^{2^n} \mid n \geq 0\}$.

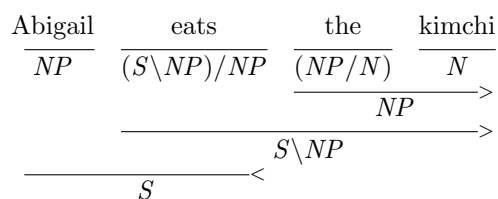


Figure 8.14: A syntactic analysis in CCG involving forward and backward function application

As with the move from regular to context-free languages, mildly context-sensitive languages can be motivated by expedience. While infinite sequences of cross-serial dependencies cannot be handled by context-free grammars, even finite sequences of cross-serial dependencies are more convenient to handle using a mildly context-sensitive formalism like **tree-adjoining grammar** (TAG) and **combinatory categorial grammar** (CCG). Furthermore, TAG-inspired parsers have been shown to be particularly effective in parsing the Penn Treebank (Collins, 1997; Carreras et al., 2008), and CCG plays a leading role in current research on semantic parsing (Zettlemoyer and Collins, 2005). Furthermore, these two formalisms are weakly equivalent: any language that can be specified in TAG can also be specified in CCG, and vice versa (Joshi et al., 1991). The remainder of the chapter gives an overview of the two formalisms, but the reader is encouraged to consult Joshi and Schabes (1997) and Steedman and Baldridge (2011) for more detail on TAG and CCG respectively.

8.3.1 Combinatory categorial grammar

In combinatory categorial grammar, structural analyses are built up through a small set of generic combinatorial operations, which apply to immediately adjacent sub-structures. These operations act on the categories of the sub-structures, producing a new structure with a new category. The basic categories include *S* (sentence), *NP* (noun phrase), *VP* (verb phrase) and *N* (noun). The goal is to label the entire span of text as a sentence, *S*.

Complex categories, or types, are constructed from the basic categories, parentheses, and forward and backward slashes: for example, *S/NP* is a complex type, indicating a sentence that is lacking a noun phrase to its right; *S\NP* is a sentence lacking a noun phrase to its left. Complex types act as functions, and the most basic combinatory operations are function application to either the right or left neighbor. For example, the type of a verb phrase, such as *eats*, would be *S\NP*. Applying this function to a subject noun phrase to its left results in an analysis of *Abigail eats* as category *S*, indicating a successful parse.

Transitive verbs must first be applied to the direct object, which in English appears on the right, before the subject, which appears on the left. They therefore have the more com-

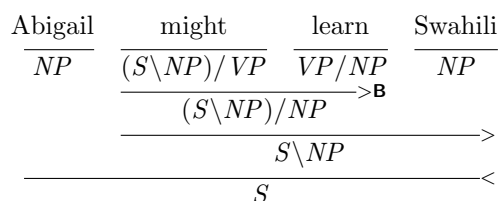


Figure 8.15: A syntactic analysis in CCG involving function composition (example based on Steedman and Baldridge, 2011)

plex type $(S \backslash NP) / NP$. Similarly, the application of a determiner to the noun at its right results in a noun phrase, so determiners have the type NP / N . An example involving a transitive verb and a determiner is shown in Figure 8.14. A key point from this example is that it can be trivially transformed into phrase-structure tree, by treating each function application as a constituent phrase. Indeed, when CCG’s only combinatory operators are forward and backward function application, it is equivalent to context-free grammar. However, the location of the “effort” has changed. Rather than designing good productions, the grammar designer must focus on the **lexicon** — choosing the right categories for each word. This makes it possible to parse a wide range of sentences using only a few generic combinatory operators.

Things become more interesting with the introduction of two additional operators: **composition** and **type-raising**. Function composition enables the combination of complex types: $X / Y \circ Y / Z \Rightarrow_B X / Z$ (forward composition) and $Y \backslash Z \circ X \backslash Y \Rightarrow_B X \backslash Z$ (backward composition).¹¹ Composition makes it possible to “look inside” complex types, and combine two adjacent units if the “input” for one is the “output” for the other. Figure 8.15 shows how function composition can be used to handle modal verbs. While this sentence can be parsed using only function application, the composition-based analysis is preferable because the unit *might learn* functions just like a transitive verb, as in the example *Abigail studies Swahili*. This in turn makes it possible to analyze conjunctions such as *Abigail studies and might learn Swahili*, attaching the direct object *Swahili* to the entire conjoined verb phrase *studies and might learn*.

Type raising converts an element of type X to a more complex type: $X \Rightarrow_T T / (T \backslash X)$ (forward type-raising to type T), and $X \Rightarrow_T T \backslash (T / X)$ (backward type-raising to type T). Type-raising makes it possible to reverse the relationship between a function and its argument — by transforming the argument into a function over functions over arguments! An example may help demystify things. Figure 8.15 shows how to analyze an object relative clause, *a story that Abigail tells*. The problem is that *tells* is a transitive verb, expecting a direct object to its right. As a result, *Abigail tells* is not a constituent in our context-free grammar of English. The issue is resolved by raising *Abigail* from NP to the complex type

¹¹The subscript **B** follows notation from Curry and Feys (1958).

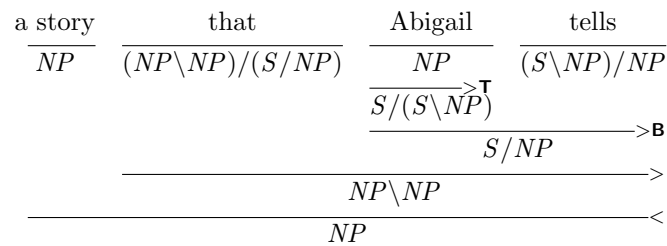


Figure 8.16: A syntactic analysis in CCG involving an object relative clause (based on slides from Alex Clark)

$(S/NP) \backslash NP$). This function can then be combined with the transitive verb *tells* by forward composition, resulting in the type (S/NP) , which is a sentence lacking a direct object to its right.¹² From here, we need only design the lexical entry for the complementizer *that* to expect a right neighbor of type (S/NP) , and the remainder of the derivation can proceed by function application.

Composition and type-raising give CCG considerable power and flexibility, but at a price. The simple sentence *Abigail tells Max* can be parsed in two different ways: by function application (first forming the verb phrase *tells Max*), and by type-raising and composition (first forming the non-constituent *Abigail tells*). This **derivational ambiguity** does not affect the resulting linguistic analysis, so it is sometimes known as **spurious ambiguity**. Hockenmaier and Steedman (2007) present a translation algorithm for converting the Penn Treebank into CCG derivations, using composition and type-raising only when necessary.

8.3.2 Tree-adjoining grammar

[todo: todo]

Exercises

1. Sketch out the state diagram for finite-state acceptors for the following languages on the alphabet $\{a, b\}$.
 - a) Even-length strings. (Be sure to include 0 as an even number.)
 - b) Strings that contain *aaa* as a substring.
 - c) Strings containing an even number of *a* and an odd number of *b* symbols.

¹²The missing direct object would be analyzed as a **trace** in CFG-like approaches to syntax, including the Penn Treebank.

- d) Strings in which the substring *bbb* must be terminal if it appears — the string need not contain *bbb*, but if it does, nothing can come after it.
2. Levenshtein edit distance is the number of insertions, substitutions, or deletions required to convert one string to another.
 - a) Define a finite-state acceptor that accepts all strings with edit distance 1 from the target string, *target*.
 - b) Now think about how to generalize your design to accept all strings with edit distance from the target string equal to d . If the target string has length ℓ , what is the minimal number of states required?
3. Construct an FSA in the style of Figure 8.3, which handles the following examples:
 - *nation*/N, *national*/ADJ, *nationalize*/V, *nationalizer*/N
 - *America*/N, *American*/ADJ, *Americanize*/V, *Americanizer*/N

Be sure that your FSA does not accept any further derivations, such as **nationalizeral* and **Americanizern*.

4. Show how to construct a trigram language model in a weighted finite-state acceptor. Make sure that you handle the edge cases at the beginning and end of the sequence accurately.
5. Extend the FST in Figure 8.6 to handle the other two parts of rule 1a of the Porter stemmer: *-sses* \rightarrow *ss*, and *-ies* \rightarrow *-i*.
6. § 8.1.4 describes T_O , a transducer that captures English orthography by transducing *cook* + *ed* \rightarrow *cooked* and *bake* + *ed* \rightarrow *baked*. Design an unweighted finite-state transducer that captures this property of English orthography.
 Next, augment the transducer to appropriately model the suffix *-s* when applied to words ending in *s*, e.g. *kiss*+*s* \rightarrow *kisses*.
7. Add parenthesization to the grammar in Figure 8.11 so that it is no longer ambiguous.
8. Construct three examples — a noun phrase, a verb phrase, and a sentence — which can be derived from the Penn Treebank grammar fragment in § 8.2.3, yet are not grammatical. Avoid reusing examples from the text. Optionally, propose corrections to the grammar to avoid generating these cases.
9. Produce parses for the following sentences, using the Penn Treebank grammar fragment from § 8.2.3.

(c) Jacob Eisenstein 2014-2017. Work in progress.

(8.21) This aggression will not stand.

(8.22) I can get you a toe.

(8.23) Sometimes you eat the bar and sometimes the bar eats you.

Then produce parses for three short sentences from a news article from this week.

10. * One advantage of CCG is its ability to handle coordination of both constituents and non-constituents:

(8.24) *Abigail and Max speak Swahili*

(8.25) *Abigail speaks and Max understands Swahili*

Define the lexical entry for *and* as

$$and := (X/X) \backslash X, \quad [8.78]$$

where X can refer to any type. Using this lexical entry, show how to parse the two examples above. In the second example, *Swahili* should be combined with the coordination *Abigail speaks and Max understands*, and not just with the verb *understands*.

Chapter 9

Context-free Parsing

Parsing is the task of determining whether a string can be produced by a given context-free grammar, and if so, how. The “how” question involves obtaining a hierarchical structure, as shown in Figure 8.13 in the previous chapter. Before we discuss specific parsing algorithms, let us consider whether exhaustive search is possible. Suppose we only have one non-terminal, X , and it has the following productions:

$$\begin{aligned} X &\rightarrow X X \\ X &\rightarrow \text{aardvark} \mid \text{abacus} \mid \dots \mid \text{zyther} \end{aligned}$$

In this grammar, the number of possible derivations for each string is equal to the number of binary bracketings, which is a **Catalan number**. Catalan numbers grow **super-exponentially** in the length of the sentence, $C_n = \frac{(2n)!}{(n+1)!n!}$. Clearly we cannot search the space of possible derivations naïvely; as with sequence labeling, we will make independence assumptions that allow us to search efficiently by reusing shared substructures with dynamic programming. This chapter will focus on a bottom-up algorithm called **CKY**; chapter 10 will describe a left-to-right algorithm called **shift-reduce**.

9.1 Deterministic bottom-up parsing

The **CKY algorithm**¹ is a bottom-up approach to parsing in a context free grammar. It efficiently tests whether a string is in a language, without considering all possible parses. The algorithm first forms small constituents, and then tries to merge them into larger constituents.

¹The name is for Cocke-Kasami-Younger, the inventors of the algorithm. It is sometimes called **chart parsing**, because of its chart-like data structure.

Let's start with a simple example grammar:

$$\begin{aligned} S &\rightarrow NP VP \\ NP &\rightarrow NP PP \mid we \mid sushi \mid chopsticks \\ PP &\rightarrow P NP \\ P &\rightarrow with \\ VP &\rightarrow V NP \mid V PP \\ V &\rightarrow eat \end{aligned}$$

Suppose we encounter the sentence *we eat sushi with chopsticks*.

- The first thing that we notice is that we can apply unary terminal productions to obtain the part-of-speech sequence NP V 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.

The CKY algorithm systematizes this approach, incrementally constructing a table t in which each cell $t[i, j]$ contains the set of nonterminals that can derive the span $w_{i:j-1}$. The algorithm fills in the upper right triangle of the table; it begins with the diagonal, which corresponds to substrings of length 1, and then computes derivations for progressively larger substrings, until reaching the upper right corner $t[0, M]$, which corresponds to the entire input. If the start symbol S is in $t[0, M]$, then the string w is in the language defined by the grammar.

- We begin by filling in the diagonal: the entries $t[m, m+1]$ for all $m \in \{0 \dots M-1\}$. These are filled with terminal productions that yield the individual tokens; for the word $w_2 = sushi$, we fill in $t[2, 3] = \{NP\}$, and so on.
- Then we fill in the next diagonal, in which each cell corresponds to a subsequence of length two: $t[0, 2], t[1, 3], \dots, t[M-2, M]$. These are filled in by looking for binary productions capable of producing at least one entry in each of the cells corresponding to left and right children. For example, the cell $t[1, 3]$ includes VP because the grammar includes the production $VP \rightarrow V NP$, and we have $V \in t[1, 2]$ and $NP \in t[2, 3]$.
- When we move to the next diagonal, there is an additional decision to make: where to split the left and right children. The cell $t[i, j]$ corresponds to the subsequence $w_{i:j-1}$, and we must choose some **split point** $i < k < j$, so that $w_{i:k-1}$ is the left child and $w_{k:j-1}$ is the right child. We do this by looping over all possible k , and

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 11 The CKY algorithm for parsing with context-free grammars

```

1: for  $m \in \{0 \dots M - 1\}$  do
2:    $t[m, m + 1] \leftarrow \{X : X \rightarrow w_m \in R\}$ 
3:   for  $n \in \{m + 1 \dots M\}$  do
4:      $t[m, n] \leftarrow \emptyset$ 
5:   for  $\ell \in \{2 \dots M\}$  do
6:     for  $m \in \{0 \dots M - \ell\}$  do
7:       for  $k \in \{m + 1 \dots m + \ell - 1\}$  do
8:          $t[m, m + \ell] \leftarrow t[m, m + \ell] \cup \{X : (X \rightarrow Y Z) \in R \wedge Y \in t[m, k] \wedge Z \in t[k, m + \ell]\}$ 
9:   if  $S \in t[0, M]$  then
10:    return True
11:  else
12:    return False

```

then looking for productions that generate elements in $t[i, k]$ and $t[k, j]$; the left-hand side of all such productions can be added to $t[i, j]$. When it is time to compute $t[i, j]$, the cells $t[i, k]$ and $t[k, j]$ have already been filled in, since these cells correspond to shorter sub-strings of the input.

- The process continues until we reach $t[0, M]$.

Algorithm 11 further formalizes this process, and Figure 9.1 shows the chart that arises from parsing the sentence *we eat sushi with chopsticks* using the grammar defined above.

An important detail about the CKY algorithm is that it assumes that all productions with non-terminals on the right-hand side (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 (see § 8.2.1). For example, if we have the production $VP \rightarrow V NP NP$ (for ditransitive verbs), we might convert to $VP \rightarrow VP_{ditrans}/NP NP$, and then add the production $VP_{ditrans}/NP \rightarrow V NP$.
- What about unary productions like $VP \rightarrow V$? In practice, this is handled by making a second pass on each diagonal, in which each cell $t[i, j]$ is augmented with all possible unary productions capable of generating each item already in the cell. Suppose our example grammar is extended to include the production $VP \rightarrow V$. Then the cell $t[1, 2]$ — corresponding to the word *eat* — would first include the set $\{V\}$, and would be augmented to the set $\{V, VP\}$ during this second pass. This would then make it possible to parse sentences like *We eat*.

Computing the parse tree We are usually interested not only in whether a sentence is in a grammar, but in what syntactic structure is revealed by parsing. As with the Viterbi

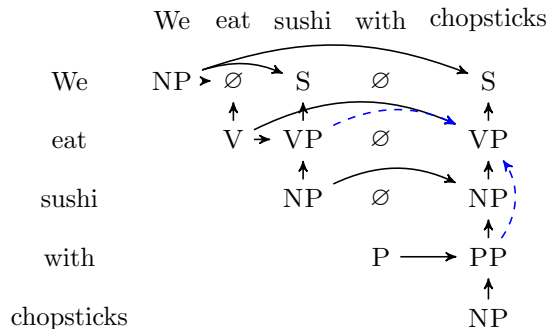


Figure 9.1: An example completed CKY chart. There are two paths to VP in position $t[1, 5]$, one in solid black and another in dashed blue.

algorithm, we can compute this structure by keeping a set of back-pointers while populating the CKY table. If we add an entry X to cell $t[i, j]$ by using the production $X \rightarrow YZ$ and the split point k , then we keep back-pointers $(i, j, X) \rightarrow (i, k, Y)$ and $(i, j, X) \rightarrow (k, j, Z)$. Once the table is constructed, we select back-pointers from $(0, M, S)$, and recursively follow them until they ground out at individual words.

For ambiguous sentences, there will be multiple paths to reach $S \in t[0, M]$. For example, in Figure 9.1, we reach $t[0, M]$ through a production that includes $VP_{int}[1, 5]$. Parsing is ambiguous for this example because there are two different ways to reach $VP \in t[1, 5]$: one with $(eat\ sushi)$ and $(with\ chopsticks)$ as children, and another with (eat) and $(sushi\ with\ chopsticks)$ as children. The presence of multiple paths indicates that the input could have been generated by the grammar in more than one way. In § 9.2, we consider different ways to resolve this ambiguity.

Complexity The space complexity of the CKY algorithm is $\mathcal{O}(M^2 \#|N|)$. We are building a table of size M^2 , and each cell must hold up to $\#|N|$ elements, where $\#|N|$ is the number of non-terminals. The time complexity is $\mathcal{O}(M^3 \#|R|)$. At each cell, we search over $\mathcal{O}(M)$ split points, and $\#|R|$ productions, where $\#|R|$ is the number of production rules in the grammar. Notice that these are considerably worse than the finite-state algorithms of Viterbi and forward-backward, which are linear time; generic shortest-path for finite-state automata has complexity $\mathcal{O}(M \log M)$.

9.2 Ambiguity in parsing

Unfortunately, parsing ambiguity is endemic to natural language:

(c) Jacob Eisenstein 2014-2017. Work in progress.

- **Attachment ambiguity:** *we eat sushi with chopsticks, I shot an elephant in my pajamas.* In each of these examples, the prepositions (*with, in*) can attach to either the verb or the direct object.
- **Modifier scope:** *southern food store, plastic cup holder.* In these examples, the first word could be modifying the subsequent adjective, or the final noun.
- **Particle versus preposition:** *The puppy tore up the staircase.* Phrasal verbs like *tore up* often include particles which could also act as prepositions.
- **Complement structure:** *The tourists objected to the guide that they couldn't hear.* This is another form of attachment ambiguity, where the complement *that they couldn't hear* could attach to the main verb (*objected*), or to the indirect object (*the guide*).
- **Coordination scope:** *"I see," said the blind man, as he picked up the hammer and saw.* In this example, the lexical ambiguity for *saw* enables it to be coordinated either with the noun *hammer* or the verb *picked up*.

These forms of ambiguity can combine, so that a seemingly simple headlines like *Fed raises interest rates* can have dozens of possible analyses, even in a minimal grammar. Broad coverage grammars permit millions of parses of typical sentences. Faced with this ambiguity, classical deterministic parsers faced a tradeoff:

- achieve broad coverage but admit a huge amount of ambiguity;
- or settle for limited coverage in exchange for constraints on ambiguity.

Rather than attempting to design a grammar that achieves broad coverage and low ambiguity, contemporary methods use labeled data to learn models capable of selecting the correct syntactic analysis.

9.2.1 Parser evaluation

Before continuing to parsing algorithms that are able to handle ambiguity, 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.² But we would like to assign *partial credit* for correctly matching parts of the reference parse. The PARSEval metrics (Grishman et al., 1992) do that, scoring each system parse via:

²Most 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.

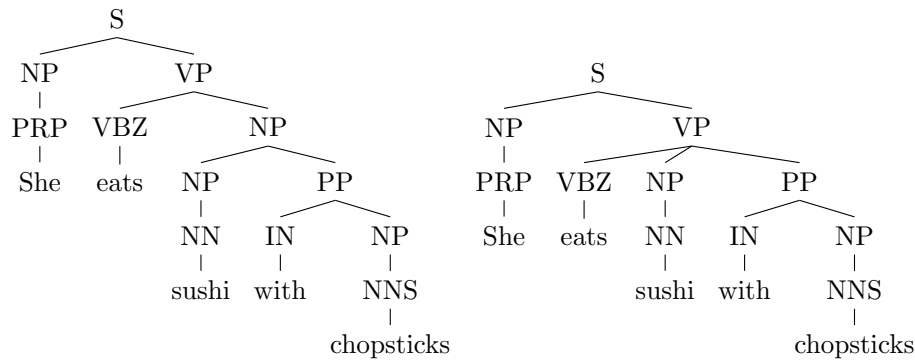


Figure 9.2: Suppose that the left parse is the system output, and the right parse is the ground truth; the precision is 0.75 and the recall is 1.0.

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.

As in chapter 3, the F-measure is the harmonic mean of precision and recall, $F = \frac{2*P*R}{R+P}$.

In **labeled** precision and recall, the system must also match the non-terminals for each bracket; in **unlabeled** precision and recall, it is only required to match the bracketing structure.

In Figure 9.2, suppose that the left tree is the system parse and the right tree is the reference parse. We have the following spans:

- $S \rightarrow w_{0:5}$ is **true positive**, because it appears in both trees.
- $VP \rightarrow w_{1:5}$ is **true positive** as well.
- $NP \rightarrow w_{2:5}$ is **false positive**, because it appears only in the system output.
- $PP \rightarrow w_{3:5}$ is **true positive**, because it appears in both trees.

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.

9.2.2 Local solutions

Some ambiguity can be resolved locally. Consider the following examples,

(9.1) [*imposed* [*a ban* [*on asbestos*]]]

(c) Jacob Eisenstein 2014-2017. Work in progress.

(9.2) [*imposed* [*a ban*] [*on asbestos*]]

This is a case of attachment ambiguity: do we attach the prepositional phrase *on asbestos* to the verb *imposed*, or the noun phrase *a ban*? To solve this problem, Hindle and Rooth (1990) proposed a likelihood ratio test:

$$LR(v, n, p) = \frac{p(p | v)}{p(p | n)} = \frac{p(on | imposed)}{p(on | ban)} \quad [9.1]$$

where they select VERB attachment if $LR(v, n, p) > 1$. The probabilities are estimated from annotated training data.

This approach is capable of modeling which prepositions tend to attach to which verbs and nouns. However, it ignores any information about the object of the prepositional phrase, which might also factor into this decision. For example:

(9.3) ...[*it* [*would end* [*its venture* [*with Maserati*]]]]

(9.4) ...[*it* [*would end* [*its venture*] [*with Maserati*]]]

The first analysis attaches *with Maserati* to the *venture*, which is almost surely preferred. Yet the likelihood ratio test ignores *Maserati*, and prefers to link the preposition *with* to the verb *end* (as in *It will end with a bang*).

- $p(with | end) = \frac{607}{5156} = 0.118$
- $p(with | venture) = \frac{155}{1442} = 0.107$

A richer probabilistic approach is undertaken by Collins and Brooks (1995), who model attachment as depending on four **heads**:

- the preposition (e.g., *with*)
- the VP attachment site (e.g., *end*)
- the NP attachment site (e.g., *venture*)
- the NP to be attached (e.g., *Maserati*)

They propose a backoff-based approach:

- First, look for counts of the tuple $\langle with, Maserati, end, venture \rangle$, and see how often VP and NP attachment were preferred for this tuple.
- If there are no counts for the full tuple, back off to the triples, $\langle with, Maserati, end \rangle + \langle with, end, venture \rangle + \langle with, Maserati, venture \rangle$, and count how often VP and NP attachment were preferred in each case.
- If there are no counts even for these triples, then try $\langle with, Maserati \rangle + \langle with, end \rangle + \langle with, venture \rangle$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Finally, if there are no counts for even these tuples, simply compute how often the preposition preferred NP or VP attachment. Since prepositions are a closed class, we can expect to have sufficient data for each preposition.

Accuracy of this method is roughly 84%. This approach of combining relative frequency estimation, smoothing, and backoff was very characteristic of 1990s statistical natural language processing. More conventional classification-based techniques can also be used for this problem; for example, Ratnaparkhi et al. (1994) designed a set of features and then trained a logistic regression classifier.

9.2.3 Beyond local solutions

Framing the problem as attachment ambiguity is limiting. It assumes the parse is mostly done, leaving just a few attachment ambiguities to solve. But realistic sentences have more than a few syntactic interpretations, and attachment decisions are interdependent. For example, consider the garden-path sentence,

(9.5) *Cats scratch people with claws with knives.*

We may want to attach *with claws* to *scratch*, as would be correct in the shorter sentence in *cats scratch people with claws*. But then we have nowhere to attach *with knives*. Only by considering these decisions jointly can we make the right choice. The task of statistical parsing is to produce a single analysis that resolves all syntactic ambiguities.

9.3 Weighted Context-Free Grammars

In a **weighted context-free grammar** (WCFG), each production $X \rightarrow \alpha$ is associated with a score $\psi_{X \rightarrow \alpha}$. The score of a derivation is simply the combination (sum or product) of the scores of all the productions. For any given string w , the “best” parse is the one corresponding to the highest-scoring derivation. More formally, for a given sequence w , we want to select the parse τ that maximizes the score,

$$\hat{\tau} = \operatorname{argmax}_{\tau: \text{yield}(\tau)=w} \sum_{(X \rightarrow \alpha) \in \tau} \psi_{X \rightarrow \alpha},$$

where we model a derivation τ as a set of productions of the form $X \rightarrow \alpha$.

As in CFGs, the **yield** of a tree is the string of terminal symbols that can be read off the leaf nodes. The set $\{\tau : w = \text{yield}(\tau)\}$ is exactly the set of all derivations of w in a CFG G .

9.3.1 Probabilistic context-free grammars

An important special case of WCFGs is a **probabilistic context-free grammar** (PCFG), in which the weight for each production $X \rightarrow \alpha$ corresponds to a log-probability, $\psi_{X \rightarrow \alpha} =$

(c) Jacob Eisenstein 2014-2017. Work in progress.

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 9.1: A fragment of an example probabilistic context-free grammar (PCFG)

$\log p(\alpha \mid X)$. These probabilities are conditioned on the left-hand side, so they must normalize to one over possible right-hand sides, $\sum_{\alpha'} p(\alpha' \mid X) = 1$. For example, for the verb phrase productions, we might have,

VP → V	0.3
VP → V NP	0.6
VP → V NP NP	0.1

which would indicate that transitive verbs are twice as common as intransitive verbs, which in turn are three times more common than ditransitive verbs.

Given probabilities on the productions, we can then score the probability of a derivation as a **product** of the probabilities of all of the productions. Consider the PCFG in

(c) Jacob Eisenstein 2014-2017. Work in progress.

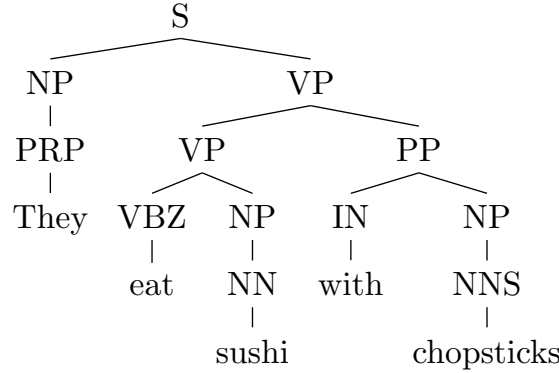


Figure 9.3: An example derivation

Table 9.1 and the parse in Figure 9.3. The probability of this parse is:³

$$\begin{aligned}
 p(\tau, \mathbf{w}) &= p(S \rightarrow NP VP) \\
 &\quad \times p(NP \rightarrow N) \times p(N \rightarrow \textit{they}) \\
 &\quad \times p(VP \rightarrow VP PP) \\
 &\quad \times p(VP \rightarrow V NP) \times p(V \rightarrow \textit{eat}) \\
 &\quad \times p(NP \rightarrow N) \times p(N \rightarrow \textit{sushi}) \\
 &\quad \times p(PP \rightarrow P NP) \times p(P \rightarrow \textit{with}) \\
 &\quad \times p(NP \rightarrow N) \times p(N \rightarrow \textit{chopsticks}) \\
 &= 0.9 \times 0.2 \times 0.2 \times 0.3 \times 0.2 \times 1.0 \times 0.2 \\
 &\quad \times \text{probability of terminal productions}
 \end{aligned}
 \tag{9.2}$$

[9.3]

Now if we consider the alternative parse in which the prepositional phrase attaches to the noun, all of these probabilities are the same, with one exception: instead of the production $VP \rightarrow VP PP$, we would have the production $NP \rightarrow NP PP$. Since $p(VP \rightarrow VP PP) > p(NP \rightarrow NP PP)$ in the PCFG, the verb phrase attachment would be preferred.⁴

9.3.2 Feature-based parsing

The scores for each production can also be computed as an inner product of weights and features,

$$\psi_{X \rightarrow \alpha} = \boldsymbol{\theta} \cdot \mathbf{f}(X, \alpha), \tag{9.4}$$

³In the remainder of the chapter, we will use the notation $p(X \rightarrow YZ)$ for the probability of producing Y and Z , conditioned on the left-hand side being X .

⁴This example hints at a big problem with PCFG parsing on non-terminals such as NP, VP, and PP: we will **always** prefer either VP or PP attachment, without regard to what is being attached! This problem is addressed later in the chapter.

where the feature vector $\mathbf{f}(X, \alpha, \mathbf{w})$ is a function of the left-hand side X and the right-hand side α . More generally, we can estimate weights for productions covering specific parts of the input,

$$\psi_{X \rightarrow \alpha, i, j, k} = \boldsymbol{\theta} \cdot \mathbf{f}(X, \alpha, \mathbf{w}, i, j, k), \quad [9.5]$$

where the feature vector is now a function of the details of the production (X and α), as well as the text \mathbf{w} and the indices of the spans to derive: the parent $\mathbf{w}_{i:j-1}$, the left child $\mathbf{w}_{i:k-1}$, and the right child $\mathbf{w}_{k:j-1}$. This is equivalent to characterizing the entire parse τ in terms of a locally-decomposable feature vector,

$$\mathbf{f}(\tau, \mathbf{w}) = \sum_{(X \rightarrow \alpha, i, j, k) \in \tau} \mathbf{f}(X, \alpha, \mathbf{w}, i, j, k) \quad [9.6]$$

$$\hat{\tau} = \operatorname{argmax}_{\tau} \boldsymbol{\theta} \cdot \mathbf{f}(\tau, \mathbf{w}) \quad [9.7]$$

$$= \operatorname{argmax}_{\tau} \sum_{(X \rightarrow \alpha, i, j, k) \in \tau} \boldsymbol{\theta} \cdot \mathbf{f}(X, \alpha, \mathbf{w}, i, j, k). \quad [9.8]$$

This enables the use of richer features, such as the words that border the span $\mathbf{w}_{i:j-1}$, the specific word at the split point w_k , the presence of a verb or noun in the left child span $\mathbf{w}_{i:j-1}$, etc. The use of such features does not affect the applicability of the CKY parsing algorithm: we can still compute each element of the table $t[i, j]$ recursively, and therefore we can still find the best parse in polynomial time. The only restriction is that the features for each production $X \rightarrow \alpha$ cannot consider other non-terminals besides the parent X and the children α . This is analogous to the Viterbi restriction that features consider only adjacent tags.

9.3.3 Estimation

Probabilistic context free grammars are similar to hidden Markov models, in that they are generative models of text. The parameters in hidden Markov models can be estimated from labeled data by relative frequency, and the same approach can be applied in PCFGs. In this case, the parameters of interest correspond to probabilities of productions, conditional on the left hand side. The relative frequency estimate is therefore,

$$p(X \rightarrow \alpha) = \frac{\text{count}(X \rightarrow \alpha)}{\text{count}(X)}. \quad [9.9]$$

For example, the probability of the production $\text{NP} \rightarrow \text{DT NN}$ is equal to the count of this production divided by the count of the non-terminal NP. This applies to terminal productions as well: the probability of $\text{NN} \rightarrow \textit{centipede}$ is the count of how often *centipede* appears in the corpus as generated from an NN tag, divided by the total count of the NN tag. These counts can be obtained from an annotated dataset, such as the **Penn Treebank**,

(c) Jacob Eisenstein 2014-2017. Work in progress.

which includes syntactic annotations over one million words of English text (Marcus et al., 1993). Even with one million words, it will be difficult to compute probabilities of relatively rare events, such as $NN \rightarrow \text{centipede}$. Therefore, smoothing techniques will again be critical for making PCFGs effective.

Feature-based parsing models can be estimated using either structure perceptron or maximum conditional likelihood. For structure perceptron, we would compute,

$$\hat{\tau} = \operatorname{argmax}_{\tau: \text{yield}(\tau) = \mathbf{w}^{(i)}} \boldsymbol{\theta} \cdot \mathbf{f}(\tau, \mathbf{w}^{(i)}) \quad [9.10]$$

$$\boldsymbol{\theta} \leftarrow \mathbf{f}(\tau^{(i)}, \mathbf{w}^{(i)}) - \mathbf{f}(\hat{\tau}, \mathbf{w}^{(i)}). \quad [9.11]$$

Alternatively, we can estimate the weights $\boldsymbol{\theta}$ by maximizing the conditional log-likelihood, $\sum_{i=1}^N \log p(\tau^{(i)} \mid \mathbf{w}^{(i)})$, which is analogous to the conditional random field (CRF) model for sequence labeling. Finkel et al. (2008) present a CRF-based parsing model. Carreras et al. (2008) present a structure perceptron model, although they perform parsing in alternative syntactic formalism known as **Tree-Adjoining Grammar** (Joshi and Schabes, 1997). § 9.5 gives more details on these discriminative parsing models.

9.3.4 Parsing with weighted context-free grammars

It is not difficult to extend the CKY algorithm to include probabilities or other weights. Let us write $\psi_{X \rightarrow Y Z}$ for the score for the production $X \rightarrow Y Z$. In the original CKY algorithm for deterministic parsing, each cell $t[i, j]$ stored a set of non-terminals capable of deriving the span $\mathbf{w}_{i:j-1}$. We now augment the table to be indexed by the tuple (i, j, X) , with $t[i, j, X]$ indicating the score of the best possible derivation of $\mathbf{w}_{i:j-1}$ from non-terminal X . Algorithm 12 shows how to perform CKY parsing using such a table.

We also keep the back-pointers corresponding to the best path to $t[i, j, X]$; the best-scoring derivation can be obtained by tracing these pointers from $t[0, M, S]$ back to each terminal, just as the best sequence of labels in the Viterbi algorithm can be computed by tracing pointers backwards from the end of the trellis. Note that we need only store back-pointers for the **best** path to $t[i, j, X]$; this follows from the locality assumption that the score for a parse is a combination of the scores of each production in the parse.

Semiring CKY As with hidden Markov models, we can generalize weighted CKY parsing using semiring notation. The basic recurrence becomes,

$$t[m, m + \ell, X] = t[m, m + \ell, X] \otimes \left(\bigoplus_{k, X \rightarrow Y Z} \psi_{X \rightarrow Y Z} \otimes t[m, k, Y] \otimes t[k, m + \ell, Z] \right). \quad [9.12]$$

This notation makes it possible to capture a number of different parsing algorithms compactly. If the scores ψ correspond to log-probabilities or feature-weight inner products,

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 12 CKY algorithm with weighted productions

```

for  $m \in \{0, \dots, M-1\}$  do
  for all  $X \in \text{tags}(w_j)$  do
     $t[m, m+1, X] \leftarrow P(X \rightarrow w_m)$ 
  for  $\ell \in \{2 \dots M\}$  do
    for  $m \in \{0, \dots, M-\ell\}$  do
      for  $k \in \{m+1, \dots, m+\ell-1\}$  do
        for all  $(X \rightarrow YZ) \in R$  do
           $t[m, m+\ell, X] \leftarrow t[m, m+\ell, X] \oplus \max_{k, X \rightarrow YZ} \psi_{X \rightarrow YZ} + t[m, k, Y] + t[k, m+\ell, Z]$ 

```

then \otimes is addition, as in Algorithm 12. If ψ are probabilities, then \otimes is multiplication. By setting \oplus equal to the max operation, the CKY algorithm computes the score of the best parse for a given sentence.

If the scores ψ correspond to probabilities, we can set \oplus to summation and \otimes to multiplication. Then the value at $t[0, M]$ corresponds to the total probability of all derivations of the input. This is known as the **inside algorithm**, and it is the tree-structured version of the **forward algorithm** from § 6.4.3. The inside algorithm can be used for unsupervised or semi-supervised parsing (Pereira and Schabes, 1992).

Finally, if we set \otimes to be the boolean “and” operation, and \oplus to be the boolean “or” operation, then $t[0, M] = \text{True}$ if and only if there is at least one derivation of the input from the grammar. Thus, the semiring notation generalizes across the weighted and unweighted CKY algorithms.

9.4 Improving Parsing by Refined Non-terminals

PCFG parsing on the Penn Treebank dataset does not perform well: Johnson (1998) shows that a PCFG estimated from treebank production counts obtains an F-measure of only $F = 0.72$. There are several problems with the use of weighted context free grammars on the Penn Treebank dataset:

- One problem is that the context-free assumption is too strict: for example, the probability of the production $\text{NP} \rightarrow \text{NP PP}$ is much higher if the parent of the noun phrase is a verb phrase (indicating that the NP is a direct object) than if the parent is a sentence (indicating that the NP is the subject of the sentence). Accurately modeling this “vertical” context is essential for accurate parsing.
- Another problem is that the Penn Treebank non-terminals are simply too coarse: there are many kinds of noun phrases and verb phrases, and accurate parsing sometimes requires knowing the difference. As we have already seen, when faced with

(c) Jacob Eisenstein 2014-2017. Work in progress.

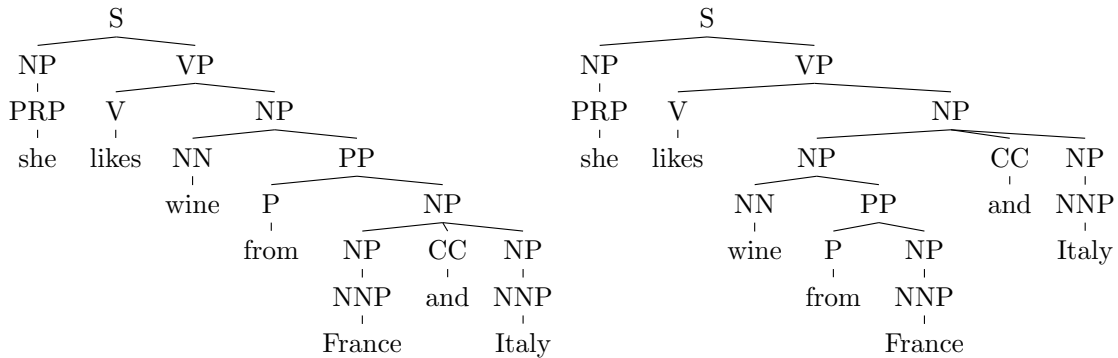


Figure 9.4: The left parse is preferable because of the conjunction of phrases headed by *France* and *Italy*.

prepositional phrase attachment ambiguity, a weighted CFG will either always choose NP attachment (if $\psi_{NP \rightarrow NP \text{ PP}} > \psi_{VP \rightarrow VP \text{ PP}}$), or it will always choose VP attachment. To get more nuanced behavior, more fine-grained non-terminals are needed.

- More generally, accurate parsing requires some amount of **semantics** — understanding the meaning of the text to be parsed. Consider the example *cats scratch people with claws*: knowledge of about *cats*, *claws*, and scratching is necessary to correctly resolve the attachment ambiguity.
- As a more extreme case, consider the example shown in Figure 9.4. The analysis on the left is preferred because of the conjunction of similar entities *France* and *Italy*. But given the non-terminals shown in the analyses, there is no way to differentiate these two parses, since they include exactly the same productions.

In all cases, what is needed seems to be more precise non-terminals. One possibility would be to rethink the linguistics behind the Penn Treebank, and ask the annotators to try again. But the original annotation effort took five years, and a more fine-grained set of non-terminals would only make things worse. We will therefore focus on automated techniques.

9.4.1 Parent annotations and other tree transformations

The key assumption underlying weighted context-free parsing is that productions depend only on the identity of the non-terminal on the left-hand side, and not on its ancestors in the parse. The validity of this assumption is an empirical question, and it depends on the non-terminals themselves: ideally, every noun phrase would be distributionally identical, so the assumption would hold. But in PTB-style analysis of English grammar, the observed probability of productions often depends on the parent of the left-hand side.

(c) Jacob Eisenstein 2014-2017. Work in progress.

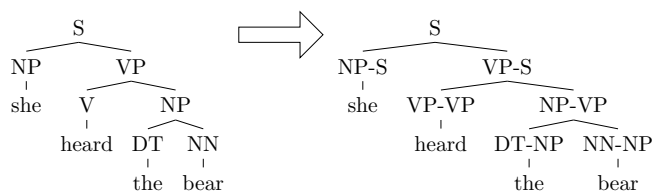


Figure 9.5: Parent annotation in a CFG derivation

For example, noun phrases are more likely to be modified by prepositional phrases when they are in the object position (e.g., *they amused the students from Georgia*) than in the subject position (e.g., *the students from Georgia amused them*). This means that the $\text{NP} \rightarrow \text{NP PP}$ production is more likely if the entire constituent is the child of a VP than if it is the child of S.

$$P(\text{NP} \rightarrow \text{NP PP}) = 11\% \quad [9.13]$$

$$P(\text{NP UNDER S} \rightarrow \text{NP PP}) = 9\% \quad [9.14]$$

$$P(\text{NP UNDER VP} \rightarrow \text{NP PP}) = 23\%. \quad [9.15]$$

Johnson (1998) proposes to capture this phenomenon via **parent annotation**. Each non-terminal is augmented with the identity of its parent, as shown in Figure 9.5). This is sometimes called **vertical Markovization**, since we introduce a Markov dependency between each node and its parent (Klein and Manning, 2003).

Using this transformation and a number of related heuristics, Johnson (1998) was able to improve the accuracy of PCFG-based parsing from 72% to 80%, at the cost of increasing the number of production rules from 14,962 to 22,773. (Recall that the number of production rules is a constant factor in the time complexity of WCFG parsing.) This increase in the number of rules is relatively modest, considering that parent annotation squares the number of non-terminals.

Parent annotation weakens the PCFG independence assumptions. This could improve accuracy by enabling the parser to make more fine-grained distinctions, which better capture real linguistic phenomena. However, each production is more rare (since the non-terminals are more specific), so the more careful smoothing is required to dampen the variance over production probabilities.

9.4.2 Lexicalization

Recall that some of the problems with PCFG parsing that were suggested above have to do with **meaning** — for example, preferring to coordinate constituents that are of the same type, like *cats* and *dogs* rather than *cats* and *houses*. A simple way to capture semantics is

(c) Jacob Eisenstein 2014-2017. Work in progress.

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 9.2: A fragment of head percolation rules

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 $\text{NP} \rightarrow \text{DT N}$ production is the noun, the head of a $\text{S} \rightarrow \text{NP VP}$ production is the head of the VP, etc. A fragment of the head percolation rules used in many parsing systems are found in Table 9.2.⁵

The meaning of the first rule 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 the modal verb *can* is tagged MD), noun phrases are headed by the rightmost noun-like non-terminal (so the head of *the red cat* is *cat*),⁶ and 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, as shown in Figure 9.6.

- In the upper part of Figure 9.6, we see how lexicalization can help solve coordination scope ambiguity. We will correctly coordinate *France* and *Italy* if,

$$p(\text{NP}(\text{Italy}) \rightarrow \text{NP}(\text{France}) \text{ CC } \text{NP}(\text{Italy})) > p(\text{NP}(\text{Italy}) \rightarrow \text{NP}(\text{wine}) \text{ CC } \text{NP}(\text{Italy})).$$

[9.16]

- In the lower part of Figure 9.6, we see how lexicalization can help solve attachment

⁵From <http://www.cs.columbia.edu/~mcollins/papers/heads>

⁶The noun phrase non-terminal is sometimes treated as a special case. Collins (1997) uses a heuristic that looks for the rightmost child which is a noun-like part-of-speech (e.g., *Nn*, *Nnp*), a possessive marker, or a superlative adjective (e.g., *the greatest*). If no such child is found, the heuristic then looks for the **leftmost** NP. If there is no child with tag NP, the heuristic then applies another priority list, this time from right to left.

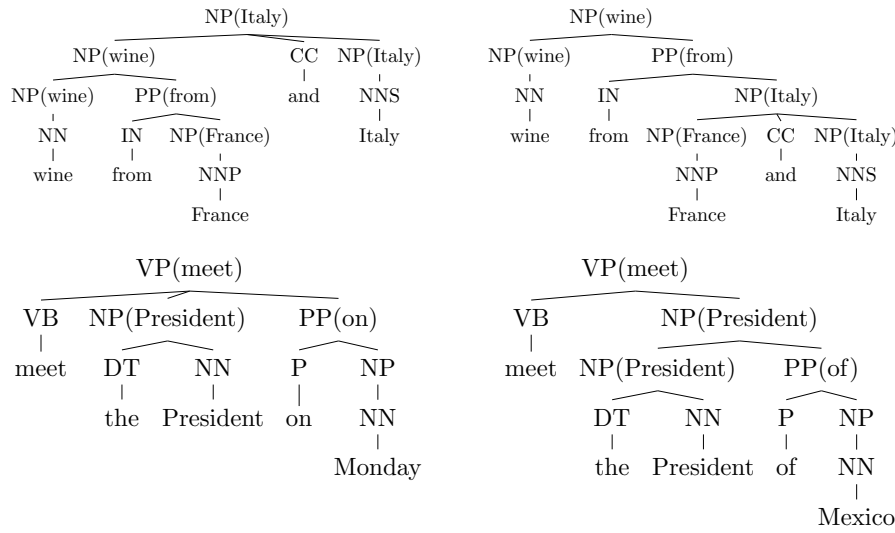


Figure 9.6: Lexicalization can address ambiguity on coordination scope (upper) and PP attachment (lower)

ambiguity. Here we assume that,

$$p(VP(meet) \rightarrow \alpha PP(on)) \gg p(NP(President) \rightarrow \beta PP(on)) \quad [9.17]$$

$$p(VP(meet) \rightarrow \alpha PP(of)) \ll p(NP(President) \rightarrow \beta PP(of)) \quad [9.18]$$

In plain English: *meetings* are usually *on* things; *Presidents* are *of* things.

- Recall that verbs may be intransitive, transitive, or ditransitive. Lexicalization can help distinguish these cases, as shown by the lexicalized PCFG probabilities for the ditransitive VP production,

$$p(VP \rightarrow V NP NP) = 0.00151 \quad [9.19]$$

$$p(VP(said) \rightarrow V(said) NP NP) = 0.00001 \quad [9.20]$$

$$p(VP(gave) \rightarrow V(gave) NP NP) = 0.01980. \quad [9.21]$$

Overall, lexicalization had a major impact on parsing accuracy, which the best lexicalized parsers attaining accuracies in the range of 87-89% (Collins, 1997, 2003; Charniak, 1997). However, lexicalized parsing introduces significant technical challenges. First, the CKY parsing algorithm must keep track of the heads of each phrase, which adds algorithmic complexity. Second, the set of possible lexicalized productions is vastly larger, since it is quadratic in the size of the vocabulary (the left and right children each have a head, and one of these heads is chosen to head the unified phrase). We now briefly overview solutions to these problems.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithms for lexicalized parsing

In weighted CFG parsing, the table element $t[i, j, X]$ keeps the score of the best derivation of the span $w_{i:j-1}$ from the non-terminal X . However, for this constituent to participate on the right-hand side of higher-level productions, we must also know its head. One solution is to expand the table to include cells of the form $t[i, j, h, X]$, where h is the index of the head of the non-terminal X for the span $w_{i:j-1}$, with $i \leq h < j$.

We can compute each element in the table by first computing the score of the best production in which the head comes from the left child, $t_\ell[i, j, h, X]$, then computing the score of the best production in which the head comes from the right child, $t_r[i, j, h, X]$, and finally taking the max over these two possibilities.

$$t_\ell[i, j, h, X] = \max_{X \rightarrow YZ} \max_{k > h} \max_{k \leq h' < j} t[i, k, h, Y] + t[k, j, h', Z] + \psi_{X(h) \rightarrow Y(h)Z(h')} \quad [9.22]$$

$$t_r[i, j, h, X] = \max_{X \rightarrow YZ} \max_{k \leq h} \max_{i \leq h' < j} t[i, k, h', Y] + t[k, j, h, Z] + \psi_{X(h) \rightarrow Y(h')Z(h)} \quad [9.23]$$

$$t[i, j, h, X] = \max(t_\ell[i, j, h, X], t_r[i, j, h, X]). \quad [9.24]$$

To compute t_ℓ , we maximize over all split points $k > h$, since the head word must be in the left child. We then maximize again over possible head words h' for the right child. An analogous computation is performed for t_r . The size of the table is now $\mathcal{O}(M^3 \#|N|)$, where M is the length of the input and $\#|N|$ is the number of non-terminals. Furthermore, each cell is computed by performing $\mathcal{O}(M^2)$ operations, since we maximize over both the split point k and the head h' . The time complexity of the algorithm is therefore $\mathcal{O}(M^5 \#|N|)$, which is impractical. Fortunately, the Eisner (1996) algorithm reduces this complexity back to $\mathcal{O}(M^3)$, using a more complex algorithm that maintains multiple tables.

The Charniak Parser

We now approach the problem of how to estimate weights for lexicalized productions $X(i) \rightarrow Y(j) Z(k)$. These productions are said to be **bilexical**, because they involve scores over pairs of words: in the example *... meet the President of Mexico*, we hope to choose the right attachment point by modeling the bilexical affinities of *(meet, of)* and *(President, of)*. The number of such word pairs is of course quadratic in the size of the vocabulary, making it difficult to estimate them directly from data.

The Charniak (1997) parser addresses this issue in the context of probabilistic parsing, so that $\psi_{X(i) \rightarrow Y(j) Z(k)}$ is equal to the (log) probability of the lexicalized production. This probability is then decomposed into a product of: a **rule probability**, which is the probability of the unlexicalized production $X \rightarrow YZ$, conditioned on the head word and parent of X (the same idea as parent annotation); a **head probability**, which is the probability of the head of X conditioned on X , the parent of X , and the head of the parent of X .

(c) Jacob Eisenstein 2014-2017. Work in progress.

Recall the example from Figure 9.6, focusing on the bottom right example, *...meet the President of Mexico*. In the case of the production $PP(of) \rightarrow P(of) NP(Mexico)$, the rule probability is $p_{\text{rule}}(PP \rightarrow P NP \mid PP, NP, of)$, since the parent is a noun phrase and the head word is *of*. The head probability is $p_{\text{head}}(of \mid PP, NP, President)$, since the parent is a noun phrase and the head of the parent is *President*. This captures the bilexical affinity between *President* and *of*, which is key to accurately parsing this example.

Even with this decomposition, it is necessary to smooth the rule and head probabilities to reduce the variance of the probability estimates. This is done by interpolating the full probabilities with simplified probabilities that condition on less information.

The Collins Parser

The Charniak parser focuses on lexical relationships between children and parents. Motivated by the linguistic theory of **lexicalized tree-adjoining grammar** (Joshi and Schabes, 1997), the Collins (2003) parser focuses on relationships between adjacent children of the same parent. We can write each production as,

$$X \rightarrow L_m L_{m-1} \dots L_1 H R_1 \dots R_{n-1} R_n,$$

where H is the child containing the head word, each L_i is a child element to the left of the head, and each R_j is a child element to the right of the head. In the Collins parser, these elements are generated probabilistically from the head outward. The outermost elements of L and R are special symbols, written \blacklozenge .

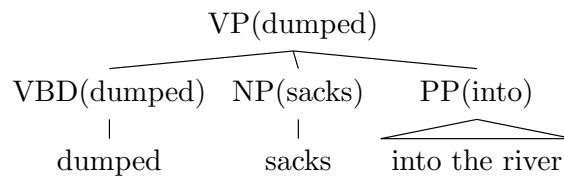


Figure 9.7: Example verb phrase for understanding the Collins parser

For example, consider the verb phrase *dumped sacks into the river*, shown in Figure 9.7. To model this rule, we would compute:

$$p(VP(dumped, VBD) \rightarrow [\blacklozenge, VBD(dumped, VBD), NP(sacks, NNS), PP(into, P), \blacklozenge]),$$

with each phrase augmented by its head word and the head word's part of speech (e.g., VBD for the head word *dumped*).

This probability is computed through a generative model, in which the head is generated first (conditioned on the parent), and then each dependent is conditioned on the

parent non-terminal and the head word. In this way, we do not directly estimate the full probability of a lexicalized production rule, but rather, we compute it from simpler probabilities involving the head and parent. Nonetheless, it is still necessary to smooth these probabilities by interpolating them with less expressive probability functions.

The Collins parser models bilexical dependencies between the head and its siblings. Bilexical probabilities require counts over pairs of words, a space of $\mathcal{O}(V^2)$ events. It is this large event space that makes these probabilities difficult to estimate, necessitating smoothing. Is it worth it? Bikel (2004) evaluates the importance of bilexical probabilities to the performance of the Collins parser. He found that bilexical probabilities are rarely available — because most of the possible bilexical pairs in the test data are unobserved in the training data — but that bilexical probabilities are indeed active in 29% of the rules in the **top-scoring** parses. Still, bilexical probabilities play a relatively small role in accuracy: an equivalent parser which conditions on only a single head suffers only 0.3% decrease in F-measure. A completely unlexicalized parser performs considerably worse, indicating that some amount of lexicalization is still necessary for top performance.

9.4.3 Refinement grammars

Lexicalization improves on pure PCFG parsing by adding detailed information in the form of lexical heads. However, estimating the probabilities of lexicalized parsing rules is difficult, requiring additional independence assumptions and complex smoothing. Klein and Manning (2003) argue that the right level of linguistic detail is somewhere between treebank categories and individual words. For example:

- 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**.
- Similarly, the words *and* and *but* should be distinguished from other coordinating conjunctions.

Figure 9.8 shows an example of an error that is corrected through the introduction of a new NP-TMP subcategory for temporal noun phrases. Klein and Manning (2003) show how the introduction of a number of such categories can make unlexicalized PCFG parsing competitive with lexicalized methods.

***Automated state-splitting** Klein and Manning (2003) use linguistic insight and error analysis to manually split PTB non-terminals so as to make parsing easier. Later work by Klein and his students automated this state-splitting process, by treating the “refined”

(c) Jacob Eisenstein 2014-2017. Work in progress.

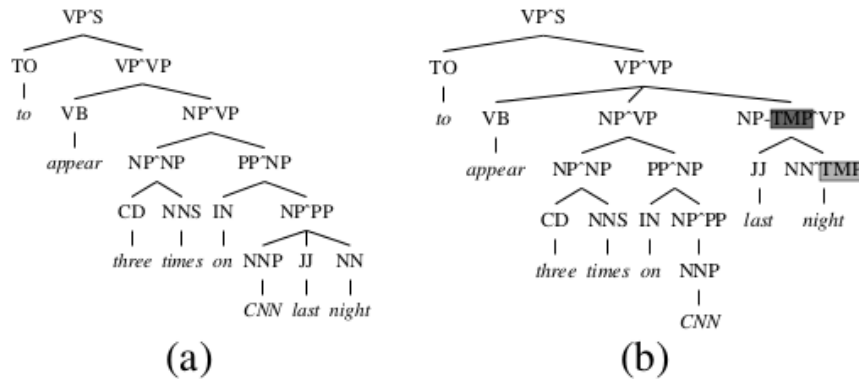


Figure 9.8: State-splitting creates a new non-terminal called NP-TMP, for temporal noun phrases. This corrects the PCFG parsing error in (a), resulting in the correct parse in (b).

non-terminals as latent variables. For example, we might split the noun phrase non-terminal into NP1, NP2, NP3, ..., without defining in advance what each refined non-terminal corresponds to.

Petrov et al. (2006) employ expectation-maximization to solve this problem. In the E-step, we estimate a marginal distribution q over the refinement type of each non-terminal. Note that this E-step is subject to the constraints of the original Penn Treebank annotation: an NP can be reannotated as NP4, but not as VP3. Now, the marginals are defined as $p(X \rightsquigarrow w_{i:j} \mid w_{1:M})$, which is the probability that the span $w_{i:j}$ is derived from the non-terminal X , conditioning on the entire sentence $w_{1:M}$ and marginalizing over all other parts of the derivation. Such marginals can be computed by a two-pass recursive algorithm called the **inside-outside algorithm** (Lari and Young, 1990).

- In the **inside** step, we compute the likelihood $p(w_{i:j} \mid X)$, which is simply the probability of deriving the span $w_{i:j}$ from the non-terminal X ; this probability depends only on the grammar, and not on any other parts of the sentence.
- In the **outside** step, we compute the probability $p(X \mid w_{1:i-1}, w_{j+1:M})$, which is the probability of a non-terminal X governing the span $w_{i:j}$, conditioned on the "outside" parts of the sentence, $w_{1:i-1}$ and $w_{j+1:M}$.

Each of these probabilities can be computed recursively. The marginal is then computed from the product of the inside and outside probabilities. The inside-outside algorithm is a direct analogue of the forward-backward algorithm, which we used to compute the marginals necessary for training conditional random fields over sequence models.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Proper nouns			
NNP-14	<i>Oct.</i>	<i>Nov.</i>	<i>Sept.</i>
NNP-12	<i>John</i>	<i>Robert</i>	<i>James</i>
NNP-2	<i>J.</i>	<i>E.</i>	<i>L.</i>
NNP-1	<i>Bush</i>	<i>Noriega</i>	<i>Peters</i>
NNP-15	<i>New</i>	<i>San</i>	<i>Wall</i>
NNP-3	<i>York</i>	<i>Francisco</i>	<i>Street</i>
Personal Pronouns			
PRP-0	<i>It</i>	<i>He</i>	<i>I</i>
PRP-1	<i>it</i>	<i>he</i>	<i>they</i>
PRP-2	<i>it</i>	<i>them</i>	<i>him</i>

Table 9.3: Examples of automatically refined non-terminals and some of the words that they generate (Petrov et al., 2006).

In the M-step, we recompute the parameters of the grammar, based on the expected counts from the E-step. As usual, this process can be iterated to convergence. To determine the number of refinement types for each tag, Petrov et al. (2006) apply a split-merge heuristic; Liang et al. (2007) and Finkel et al. (2007) apply Bayesian nonparametrics.

This approach yielded state-of-the-art accuracy at the time. Some examples of refined non-terminals are shown in Table 9.3. The proper nouns differentiate months, first names, middle initials, last names, first names of places, and second names of places; each of these will tend to appear in different parts of grammatical productions. The personal pronouns differentiate grammatical role, with PRP-0 appearing in subject position at the beginning of the sentence (note the capitalization), PRP-1 appearing in subject position but not at the beginning of the sentence, and PRP-2 appearing in object position.

9.5 Discriminative parsing

The methods described in the previous section are all based on generative parsing models, in which the probability of a parse is a product of the probabilities of the individual productions. As we have seen, these models can be improved by using finer-grained non-terminals, via parent-annotation, lexicalization, and state-splitting. An alternative path to making parsing more accurate is to use techniques from discriminative machine learning. With the exception of reranking (discussed below), the introduction of discriminative methods to parsing came relatively late. The main reason is that these learning algorithms require multiple passes over the data, applying the parser repeatedly. Unlike sequence labeling, where the time complexity of inference is linear in the size of the input, the cost of inference for parsing is non-trivial — cubic in the length of the input. These limitations

(c) Jacob Eisenstein 2014-2017. Work in progress.

prevented well-known discriminative learning techniques, such as structured perceptron, from being applied sooner.

9.5.1 Reranking

An inexpensive way to get the benefits of discriminative learning is through **reranking** (Charniak and Johnson, 2005; Collins and Koo, 2005). First, a generative model — such as the Collins or Charniak parser — is used to identify the K -best parses for a sentence. (A modified version of CKY can compute the K -best parses efficiently.) Then a discriminative learning algorithm is trained to select the best of these parses. The discriminative model does not need to search over all parses, it only needs to consider the best K identified by the “generator.” This means that the discriminator 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. Because learning is discriminative, rerankers can also use rich lexicalized features, relying on regularization to combat overfitting. Overall, this approach yields substantial improvements in accuracy on the Penn Treebank, and can be applied to improve any generative parsing model. The main limitation is that reranking can only find the best parse among the K -best offered by the generator, so it is inherently limited by the ability of the generator to find high-quality parse candidates.

9.5.2 Parsing as structure prediction

As shown in § 9.3.4, the weights on productions need not correspond to probabilities; the CKY algorithm can apply to **any** set of weights, as long as they are context-free. Discriminative learning can therefore be applied by setting $\psi_{X \rightarrow YZ} = \theta \cdot f(X \rightarrow YZ, w, i, j, k)$, with the indices i, j, k indicating the boundaries of the parent ($w_{i:j-1}$) and its left and right children ($w_{i:k-1}$ and $w_{k:j-1}$). Such features could incorporate lexical information, so that we learn weights for non-terminal productions as well as for lexicalized forms. For example:

- $f1: NP(*) \rightarrow NP(*) PP(*)$
- $f2: NP(cats) \rightarrow NP(cats) PP(*)$
- $f3: NP(*) \rightarrow NP(*) PP(claws)$
- $f4: NP(cats) \rightarrow NP(cats) PP(claws)$

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 weighted CFG 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. Like CKY, the runtime of the inside-outside algorithm is cubic in the length of the input. Because each instance must be visited and parsed many times during stochastic gradient descent, efficiency is critical. One solution is to “prefilter” the CKY parsing chart, identifying and eliminating 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; Carreras et al., 2008). They use features that capture “grandparent” dependencies between words and the heads of their parents’ parents. These second-order dependency features make the time complexity $\mathcal{O}(M^4)$ in the length of the input, so pruning is again required to make parsing efficient enough to train accurately.

9.5.3 Neural parsing

Recent work has applied neural representations to parsing, representing units of text with dense numerical vectors (Socher et al., 2013a; Durrett and Klein, 2015).

Chapter 10

Dependency Parsing

The previous chapter discussed algorithms for analyzing sentences in terms of nested **constituents**, such as noun phrases and verb phrases. The combination of constituency structure and head-percolation rules yields a set of **dependencies** between individual words. These dependencies are a more “bare-bones” version of syntax, leaving out information that is present in the full constituent parse. Nonetheless, the dependency representation is still capable of capturing important linguistic phenomena, such as the prepositional phrase attachment and coordination scope. For this reason, dependency parsing is increasingly used in applications that require syntactic analysis. While dependency structures can be obtained as a byproduct of constituent parsing, it is more efficient to extract them directly. Indeed, accurate dependency parses can be obtained by algorithms with time complexity that is linear in the length of the sentence. This chapter begins by overviewing dependency grammar, and then presents the two dominant approaches to dependency parsing, graph-based and transition-based dependency parsing.

10.1 Dependency grammar

In lexicalized parsing, non-terminals such as NP are augmented with **head words**, as shown in Figure 10.1a. In this sentence, the head of the S constituent is the main verb, *scratch*; this non-terminal then produces the noun phrase *the cats*, whose head word is *cats*, and from which we finally derive the word *the*. Thus, the word *scratch* occupies the central position for the sentence, with the word *cats* playing a supporting role. In turn, *cats* occupies the central position for the noun phrase, with the word *the* playing a supporting role.

These relationships, which hold between the words in the sentence, can be formalized in a directed graph structure. In this graph, there is an edge from word i to word j iff word i is the head of the first branching node above a node headed by j . Thus, in our

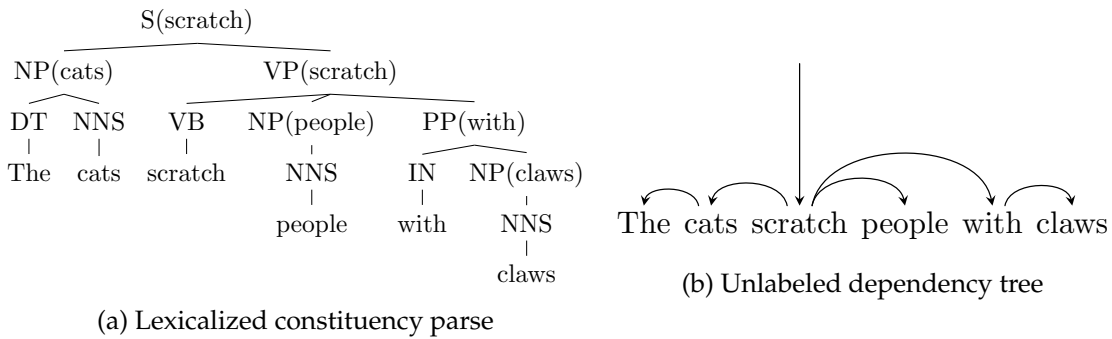


Figure 10.1: Dependency grammar is closely linked to lexicalized context free grammars: each lexical head has a dependency path to every other word in the constituent.

example, we would have *scratch* \rightarrow *cats* and *cats* \rightarrow *the*. We would not have the edge *scratch* \rightarrow *the*, because although *scratch* dominates *the* in the graph, it is not the head of a node that produces a node headed by *the*. These edges describe syntactic **dependencies**, a bilexical relationship between a **head** and a **dependent**, which is at the heart of **dependency grammar** (Tesnière, 1966).

If we continue to build out this **dependency graph**, we will eventually reach every word in the sentence, as shown in Figure 10.1b. In this graph — and in all graphs constructed in this way — every word will have exactly one incoming edge, except for the root word, which is indicated by a special incoming arrow from above. Another feature of this graph is that it is **weakly connected**, in the sense that if we replaced the directed edges with undirected edges, there would be a path between all pairs of nodes. From these properties, it can be shown that there are no cycles in the graph (or else at least one node would have to have more than one incoming edge), and therefore, the graph is a **tree**.

Although we have begun by motivating dependency grammar in terms of lexicalized constituent parsing, there is a rich literature on dependency grammar as a model of syntax in its own right (Tesnière, 1966). Kübler et al. (2009) provides a comprehensive overview of this literature.

10.1.1 What do the edges mean?

A dependency edge implies an asymmetric syntactic relationship between the head and dependent words. For a pair like *the cats* or *cats scratch*, how do we decide which is the head? Here are some possible criteria:

- The head sets the syntactic category of the construction: for example, nouns are the heads of noun phrases, and verbs are the heads of verb phrases.

(c) Jacob Eisenstein 2014-2017. Work in progress.

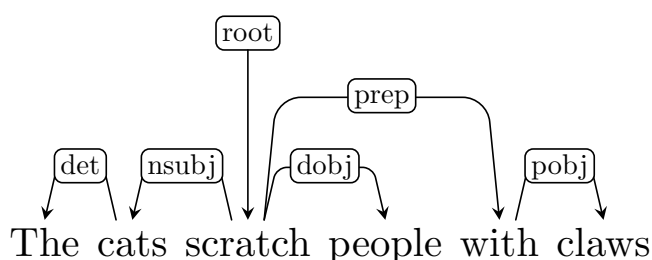


Figure 10.2: A labeled dependency parse

- The modifier may be optional while the head is mandatory: for example, in the sentence *cats scratch people with claws*, the substrings *cats scratch* and *cats scratch people* are grammatical sentences, but *with claws* is not.
- The head determines the morphological form of the modifier: for example, in languages that require gender agreement, the gender of the noun determines the gender of the adjectives and determiners.

As always, these guidelines sometimes conflict, but it is possible to use these basic principles to define fairly consistent conventions at the level of part-of-speech tags, similar to the head percolation rules from lexicalized constituent parsing.

Edges may be **labeled** to indicate the nature of the syntactic relation that holds between the two elements. An example is shown in Figure 10.2. The edge between *scratch* and *cats* is labeled NSUBJ, with *scratch* as the head; this indicates that the noun subject of the predicate verb *scratch* is headed by the word *cats*. The edge from *scratch* to *people* is labeled with DOBJ; this indicates that the word *people* is the head of the direct object. The Stanford typed dependencies have become a standard inventory of dependency types for English (De Marneffe and Manning, 2008). De Marneffe et al. (2014) propose a more minimal “universal” set of dependencies that is suitable for many languages.

10.1.2 Ambiguity and difficult cases

[**todo: update this section with current standards from universal dependency treebank (Nivre et al., 2016)**] The attachment ambiguity in the sentence shown in Figure 10.2 can be represented by a single change: replacing the edge from *scratch* to *with* by an edge from *people* to *with*. This should give you an idea of why labeled dependency trees are useful: they tell us who did what to whom.

However, dependency trees are less structurally expressive than lexicalized CFG derivations. That means they hide information that would be present in a CFG parse. Often this “information” is in fact irrelevant for any conceivable linguistic purpose: for example, Figure 10.3 shows three different ways of representing prepositional phrase adjuncts to

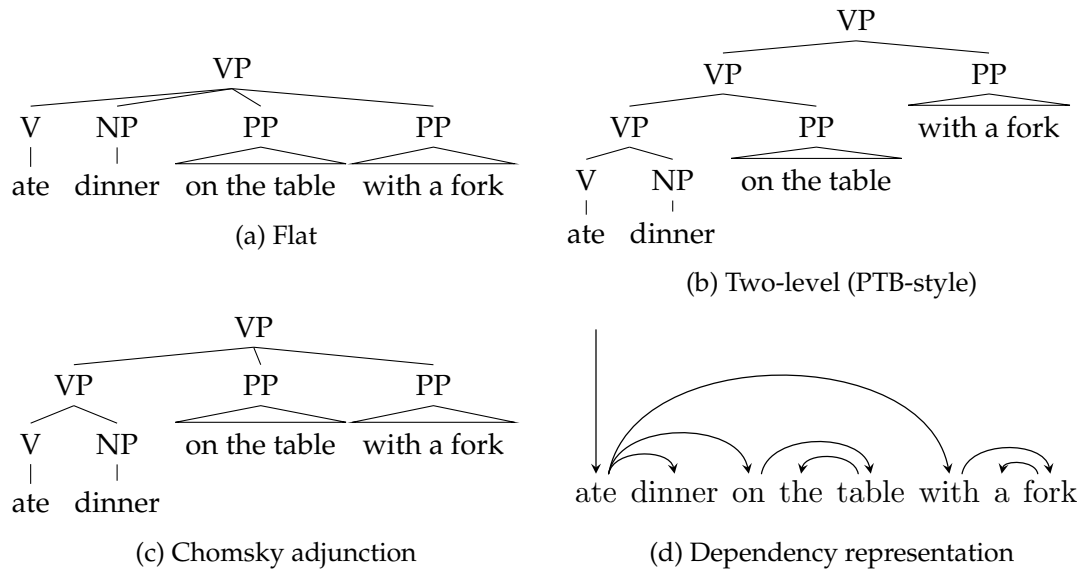


Figure 10.3: The three different CFG analyses of this verb phrase all correspond to a single dependency structure.

the verb *ate*. Because there is apparently no meaningful difference between these analyses, the Penn Treebank decides by convention to use the two-level representation. As shown in Figure 10.3d, these three cases all look the same in a dependency parse. So if you didn't think there was any meaningful difference between these three constituent representations, you may view this as an advantage of the dependency representation.

Dependency grammar still leaves open some tricky representational decisions. For example, coordination is a challenge: in the sentence, *Abigail and Max like kimchi* (Figure 10.4), which word is the immediate dependent of the main verb *likes*? Choosing either *Abigail* or *Max* seems arbitrary; for fairness we might choose *and*, but this seems in some ways to be the least important word in the noun phrase. One typical solution is to simply choose the left-most item in the coordinated structure — in this case, *Abigail*. Another alternative, as shown in Figure 10.4c, is a **collapsed** dependency grammar in which conjunctions are not included as nodes in the graph, but are instead used to label the edges (De Marneffe et al., 2006). Popel et al. (2013) survey alternatives for handling this phenomenon across several dependency treebanks.

The same logic that makes us reluctant to accept *and* as the head of a coordinated noun phrase may also make us reluctant to accept a preposition as the head of a prepositional phrase. In the sentence *cats scratch people with claws*, surely the word *claws* is more central than the word *with* — and it is precisely the bilocal relations between *scratch*, *claws*, and *people* that help guide us to the correct syntactic interpretation. Yet there are also argu-

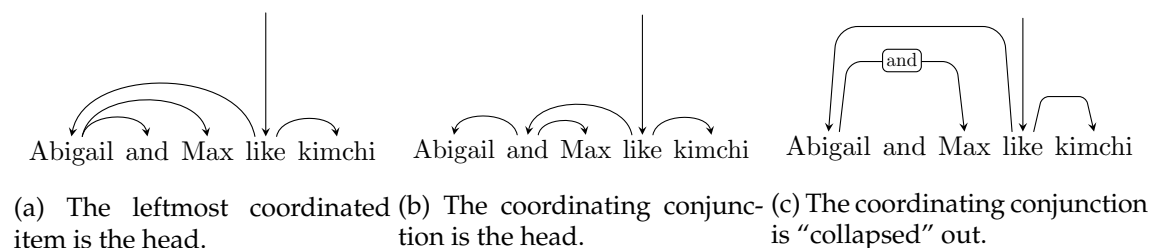


Figure 10.4: Three alternatives for representing coordination in a dependency parse

ments for preferring the preposition as the head — as we saw in § 9.4.2, the preposition itself is what helps us to choose verb attachment in *meet the President on Monday* and noun attachment in *meet the President of Mexico*. Collapsed dependency grammar is again a possible solution: we can collapse out the prepositions so that the dependency chain,

$$President \rightarrow_{prep} of \rightarrow_{pobj} Mexico$$

would be replaced by $President \rightarrow_{PREP:of} Mexico$. Dependency annotation is an active area of research due to the ongoing development of the universal dependency treebank, which has produced dependency-annotated corpora in 47 languages at the time of this writing (Nivre et al., 2016).¹

10.1.3 Projectivity

The dependency graphs that can be built from all possible lexicalized constituent parses of a sentence with M words are a proper subset of the spanning trees over M nodes. In other words, there exist spanning trees that do not correspond to any lexicalized constituent parse. This is because syntactic constituents are **contiguous** spans of text, so that the head h of the constituent that spans the nodes from i to j must have a path to every node in this span. This property is known as **projectivity**. Informally, it means that “crossing edges” are prohibited. The formal definition follows:

Definition 2 (Projectivity). *An edge from i to j is projective iff all k between i and j are descendants of i . A dependency parse is projective iff all its edges are projective.*

If we were to annotate a dependency parse directly — rather than deriving it from a lexicalized constituent parse — such non-projective edges would occur. Figure 10.5 gives an example of a non-projective dependency graph in English. This dependency graph does not correspond to any constituent parse. In languages where non-projectivity is common, such as Czech and German, it is better to annotate dependency trees directly,

¹<http://universaldependencies.org/>

	% non-projective edges	% non-projective sentences
Czech	1.86%	22.42%
English	0.39%	7.63%
German	2.33%	28.19%

Table 10.1: Frequency of non-projective dependencies in three languages (Kuhlmann and Nivre, 2010)

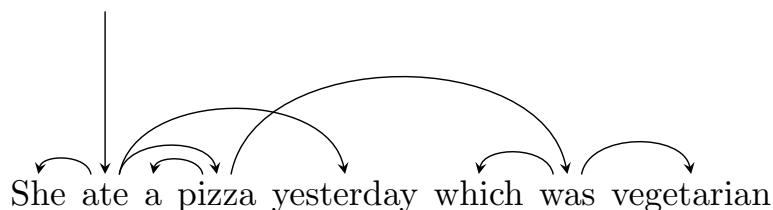


Figure 10.5: An example of a non-projective dependency parse in English

rather than deriving them from constituent parses. An example is the Prague Dependency Treebank (Böhmová et al., 2003), which contains 1.5 million words of Czech, with approximately 12,000 non-projective edges (see Table 10.1). Even though relatively few dependencies are non-projective in Czech and German, many sentences have at least one such dependency.

As we will see in the next section, projectivity has important consequences for the sorts of algorithms that can perform dependency parsing.

10.2 Graph-based dependency parsing

Let $\mathbf{y} = \{\langle i, j, r \rangle\}$ indicate a dependency graph with relation r from head word w_i to dependent word w_j . We would like to define a scoring function $\boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{y}, \mathbf{w})$, where $\mathbf{f}(\mathbf{y}, \mathbf{w})$ is a vector of features on the dependency graph and sentence, and $\boldsymbol{\theta}$ is a vector of weights. The dependency parsing problem is then the structure prediction problem,

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathcal{Y}(\mathbf{w})} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{y}, \mathbf{w}). \quad [10.1]$$

As usual, the number of possible labelings $\mathcal{Y}(\mathbf{w})$ is exponential in the length of the input. In the case of non-projective dependency parsing, the set $\mathcal{Y}(\mathbf{w})$ includes all possible spanning trees over a complete graph with M nodes, where M is the length of the sentence \mathbf{w} . The size of this set is M^{M-2} (Wu and Chao, 2004). Algorithms that search over this space of possible graphs are known as **graph-based dependency parsers**.

(c) Jacob Eisenstein 2014-2017. Work in progress.

In sequence labeling and constituent parsing, it was possible to search efficiently over an exponential space by choosing a feature function that decomposes into a sum of local feature vectors. A similar approach is possible for dependency parsing, by requiring the feature function to decompose across dependency arcs $i \rightarrow j$:

$$f(y, w) = \sum_{\langle i, j, r \rangle \in y} f(w, i, j, r) \quad [10.2]$$

$$\theta \cdot f(y, w) = \sum_{\langle i, j, r \rangle \in y} \theta \cdot f(w, i, j, r). \quad [10.3]$$

Dependency parsers that operate under this assumption are known as **arc-factored**, since the overall (exponentiated) score is a product of scores over all arcs. As described later in this section, the arc-factored assumption enables efficient algorithms for dependency parsing.

10.2.1 Features

Typical features for arc-factored dependency parsing are similar to those used in sequence labeling and discriminative constituent parsing. They include: the length and direction of the dependency arc; the words linked by the dependency relation; their prefixes, suffixes, and part-of-speech tags (as produced by an automatic tagger); and their neighbors in the sentence. In labeled dependency parsing, each of these features are also conjoined with the relation type r .

Bilexical features, which include both the head and the dependent, will be helpful for common words, but will be extremely sparse for rare words. It is therefore necessary to include features at various levels of detail, such as: word-word, word-tag, tag-word, and tag-tag. For example, for the arc *scratch* \rightarrow *cats*, we might have the features,

$$\begin{array}{ll} \{w_i \rightarrow w_j : & \text{scratch} \rightarrow \text{cats}, \\ w_i \rightarrow t_j : & \text{scratch} \rightarrow \text{NNS}, \\ t_i \rightarrow w_j : & \text{VBP} \rightarrow \text{cats}, \\ t_i \rightarrow t_j : & \text{VBP} \rightarrow \text{NNS}\} \end{array}$$

Regularized discriminative learning algorithms can then learn to trade off between features that are rare but highly predictive, and features that are common but less informative.

As with sequence labeling, it is possible to include features on neighboring words without breaking the locality restriction: we can consider features such as the identity, part-of-speech, and shape of the preceding and succeeding words, $w_{i-1}, w_{i+1}, w_{j-1}, w_{j+1}$. What we cannot do (yet) is consider other parts of the graph y , such as the parent of i

(which I will denote $w_{\Gamma(i)}$) or the siblings of j , the set $\{w_j : \Gamma(j) = i\}$. This requires higher-order dependency parsing, discussed in § 10.2.5.

To give a concrete example, the seminal paper by McDonald et al. (2005a) includes the following features for an arc between words w_i and w_j , with part-of-speech tags t_i and t_j :

Unigram features $\langle w_i \rangle; \langle t_i \rangle; \langle w_i, t_i \rangle; \langle w_j \rangle; \langle t_j \rangle; \langle w_j, t_j \rangle$.

Bigram features $\langle w_i, t_i, w_j, t_j \rangle; \langle w_i, w_j, t_j \rangle; \langle t_i, w_j, t_j \rangle; \langle w_i, t_i, t_j \rangle; \langle w_i, t_i, w_j \rangle; \langle w_i, w_j \rangle; \langle t_i, t_j \rangle$.

“In-between” features $\langle t_i, t_k, t_j \rangle$ for all k between i and j .

Neighbor features $\langle t_i, t_{i+1}, t_{j-1}, t_j \rangle; \langle t_{i-1}, t_i, t_{j-1}, t_j \rangle; \langle t_i, t_{i+1}, t_j, t_{j+1} \rangle; \langle t_{i-1}, t_i, t_j, t_{j+1} \rangle$

In addition, all the word features are supplemented with the five-character prefixes for all words longer than five characters (e.g., *unconscionable* \rightarrow *uncon*). The bigram features include several varieties of backoff from the most detailed 4-tuple feature; McDonald et al. (2005a) note that these backoff features were particularly helpful, presumably because they improve generalization. The “in-between” features activate for all part-of-speech tags between positions i and j in the sentence. This feature group helps to “rule out situations when a noun would attach to another noun with a verb in between, which is a very uncommon phenomenon.”

10.2.2 Learning

Having formulated graph-based dependency parsing as a structure prediction problem, we can apply similar learning algorithms to those used in sequence labeling. The most direct application is structured perceptron,

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y}' \in \mathcal{Y}(\mathbf{w})} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, \mathbf{y}') \quad [10.4]$$

$$\boldsymbol{\theta} = \boldsymbol{\theta} + \mathbf{f}(\mathbf{w}, \mathbf{y}) - \mathbf{f}(\mathbf{w}, \hat{\mathbf{y}}) \quad [10.5]$$

This is just like sequence labeling, but now the argmax requires a maximization over all dependency trees for the sentence. Algorithms for performing this search efficiently are described below. We can apply all the usual tricks from § 1.2: weight averaging, large-margin, and regularization. McDonald et al. (2005a,b) were the first to treat dependency parsing as a structure prediction problem, using MIRA (a close relative of the passive-aggressive algorithm we saw in § 1.2) to obtain high accuracy parses in both projective and non-projective settings.

Conditional random fields (CRFs) are globally-normalized conditional models (see chapter 6), and they can be applied to any graphical model in which we can efficiently compute marginal probabilities over individual random variables — in this case, we need

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 13 Chu-Liu-Edmonds algorithm for unlabeled dependency parsing

```

1: procedure CHU-LIU-EDMONDS( $\{\psi(i \rightarrow j)\}_{i,j \in \{1 \dots M\}}$ )
2:   for  $j \in 1 \dots M$  do
3:      $h_j \leftarrow \operatorname{argmax}_i \psi(i \rightarrow j)$  ▷ Find the best incoming edge for each node
4:    $\tau \leftarrow \{j, h_j\}_{j \in 1 \dots M}$  ▷  $\tau$  is the graph of the best incoming edges
5:    $\mathcal{C} \leftarrow \text{FINDCYCLES}(\tau)$ 
6:   if  $\mathcal{C} = \emptyset$  then
7:     return  $\tau$  ▷ If  $\tau$  has no cycles, it is the best tree
8:   else ▷ Otherwise, collapse each cycle
9:     for each cycle  $c \in \mathcal{C}$  do
10:       Remove all nodes in the cycle from the graph
11:       Add a “super-node” representing the cycle
12:     Let  $G$  be the resulting graph
13:   return CHU-LIU-EDMONDS( $G$ ) ▷ Call recursively on the collapsed graph

```

marginals over the edges. The marginals are required because the unregularized log-likelihood has a gradient that sums over all possible edges, taking the difference between the features in the observed dependency parses and the expected feature counts under $p(\mathbf{y} \mid \mathbf{w})$:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{\theta}} = \sum_{(i,j) \in \mathcal{V}} \mathbf{f}(\mathbf{w}, i, j) - \sum_{i,j} p(i \rightarrow j \mid \mathbf{w}) \mathbf{f}(\mathbf{w}, i, j) \quad [10.6]$$

For projective dependency trees, the marginal probabilities can be computed in cubic time, using a variant of the inside-outside algorithm (Lari and Young, 1990). For non-projective dependency parsing, marginals can also be computed in cubic time, using the **matrix-tree theorem** (Koo et al., 2007; McDonald et al., 2007; Smith and Smith, 2007). Details of these methods are described by Kübler et al. (2009).

10.2.3 Algorithms for non-projective dependency parsing

In **non-projective dependency parsing**, the goal is to identify the highest-scoring spanning tree over the words in the sentence. The arc-factored assumption ensures that the score for each spanning tree will be computed as a sum over scores for the edges. We can precompute these scores, $\psi(i \rightarrow j, r) = \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, i, j, r)$, before applying a parsing algorithm. (We must compute $\mathcal{O}(M^2 R)$ such scores, where M is the length of the sentence and R is the number of dependency relation types, so this is a lower bound on the time complexity of any exact algorithm for dependency parsing.)

Based on these scores, we build a weighted connected graph. Arc-factored non-projective dependency parsing is then equivalent to finding the the spanning tree that achieves the

(c) Jacob Eisenstein 2014-2017. Work in progress.

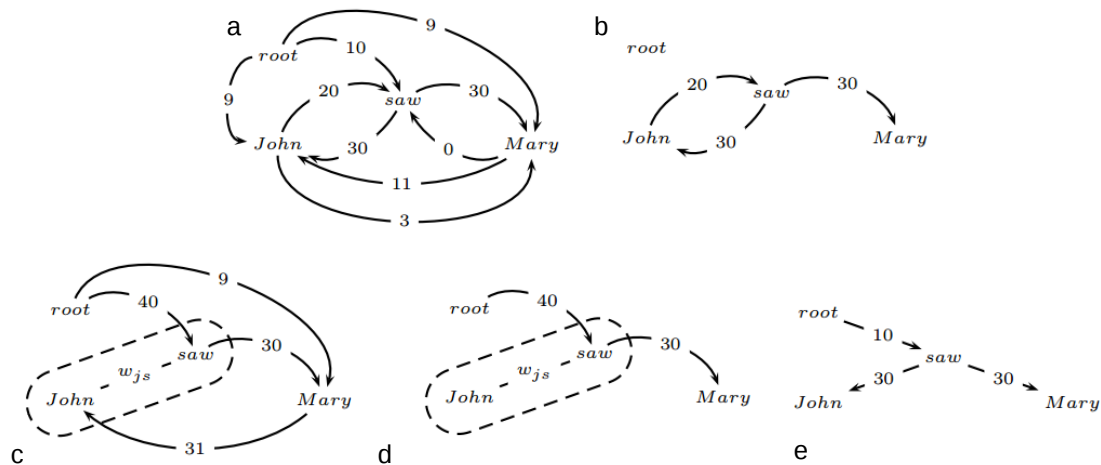


Figure 10.6: An illustration of the MST algorithm on a simple example. Figure borrowed from McDonald et al. (2005b).

maximum total score, $\sum_{\langle i,j,r \rangle \in \mathbf{y}} \psi(i \rightarrow j, r)$. The **Chu-Liu-Edmonds algorithm** (Chu and Liu, 1965; Edmonds, 1967) computes this spanning tree in time $\mathcal{O}(M^3 R)$. The algorithm, which is sketched in Algorithm 13, operates recursively. It first identifies the highest scoring incoming edge for each node, and then checks the graph for cycles. If there are no cycles, the resulting graph is a spanning tree, and moreover, it is the maximum spanning tree, because there is no better-scoring incoming edge for each node. If there is a cycle, the cycle is collapsed into a “super-node”, whose incoming edges have scores equal to the scores of the best spanning tree that includes both the edge and all nodes in the cycle. [todo: more detail on what happens when you collapse cycles].

The algorithm works because it can be proved that the maximum spanning tree on the contracted graph is equivalent to the maximum spanning tree on the original graph. The basic process is illustrated in Figure 10.6. In part (a), we see the complete graph, which

(c) Jacob Eisenstein 2014-2017. Work in progress.

includes all edge scores $\psi(i \rightarrow j)$. In (b), we see the highest scoring incoming edge for each node. In (c), the cycle between *John* and *saw* is contracted, creating new incoming edges with weight 40 from the root, and weight 31 from *Mary*. In (d), we find the highest-scoring incoming edge in the new graph. There are no remaining cycles, so we recover the maximum spanning tree.

Let us consider the time complexity of unlabeled dependency parsing first. For each of the M words in the sentence, one must search all $M - 1$ other words for the highest-scoring incoming edge, for a time complexity of $\mathcal{O}(M^2)$. In the worst case, it is necessary to contract the graph M times. If we redo the search within each contraction, we face a total cost of $\mathcal{O}(M^3)$. Recall that the CKY constituent parsing algorithm is also cubic time complexity in the length of the sentence. However, further optimizations are possible, resulting in a complexity of $\mathcal{O}(M^2)$ (Tarjan, 1977). To generalize the algorithm to labeled dependency parsing, it is necessary only to compute the best scoring label for each possible edge. Because of the arc-factoring assumption, the edge labels are decoupled from each other, so this can be done as a preprocessing step, with total complexity of $\mathcal{O}(M^2 R)$.

10.2.4 Algorithms for projective dependency parsing

The Chu-Liu-Edmonds algorithm finds the best scoring dependency tree, but it does not enforce the projectivity constraint. For languages in which we expect projectivity — such as English — we may prefer to ensure that the parsing algorithm returns only projective trees. Note that the arc-factored assumption makes it impossible to **learn** to produce projective trees, since projectivity cannot be encoded in a feature that decomposes over individual arcs.

Recall that it is possible to convert any lexicalized constituent parse directly into a projective dependency parse. This means that any algorithm for lexicalized constituent parsing is also an algorithm for projective dependency parsing. One such algorithm is presented in § 9.4, in which we built a table where the cell $t[i, j, h, X]$ contains the score of the best derivation of the substring $w_{i:j}$ from non-terminal X , in which the head is w_h . For unlabeled projective dependency parsing, we can apply a very similar algorithm:

$$t_\ell[i, j, h] = \max_{k > h} \max_{k \leq h' < j} t[i, k, h] + t[k, j, h'] + \psi(h \rightarrow h') \quad [10.7]$$

$$t_r[i, j, h] = \max_{k \leq h} \max_{i \leq h' < k} t[i, k, h'] + t[k, j, h] + \psi(h \rightarrow h') \quad [10.8]$$

$$t[i, j, h] = \max(t_\ell[i, j, h], t_r[i, j, h]). \quad [10.9]$$

The goal is for $t[i, j, h]$ to contain the score of the best-scoring projective dependency tree for $w_{i:j}$, headed by w_h . We must first maximize over all h' , which is the location of an immediate dependent of w_h . Projectivity guarantees that the subtree headed by h' will extend to one of the endpoints of the entire span: either from the left endpoint i to some

midpoint k , or from some midpoint k to the right endpoint j . We compute the best score for each of these possibilities separately in Equation 10.7 and Equation 10.8. Computing each of these scores also involves maximizing over all possible midpoints k .

We construct the table t from the bottom up: first compute scores for all subtrees of size 2, then size 3, and so on. The total size of the table is $\mathcal{O}(M^3)$, and to complete each cell we must search over $\mathcal{O}(M)$ dependents and $\mathcal{O}(M)$ split points. Thus, the overall complexity is $\mathcal{O}(M^5)$. The Eisner (1996) algorithm reduces this complexity to $\mathcal{O}(M^3)$ by maintaining multiple tables. For a detailed description of this algorithm, see Kübler et al. (2009). As with the Chu-Liu-Edmonds algorithm, the best-scoring label for each edge can be computed as a preprocessing step. [todo: write up formal algorithm description]

10.2.5 Higher-order dependency parsing

Arc-factored dependency parsers can only score dependency graphs as a product across their edges. However, it can be useful to consider higher-order features, which consider pairs or triples of edges, as shown in Figure 10.7. Second-order features consider **siblings** and **grandchildren**; third-order features consider **grand-siblings** (siblings and grandparents together) and **tri-siblings**.

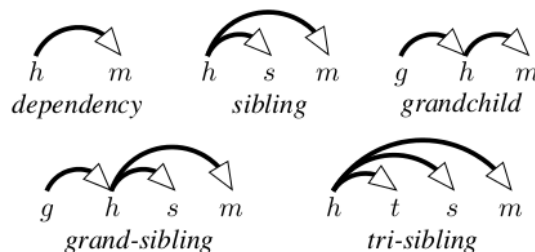


Figure 10.7: Feature templates for higher-order dependency parsing (Koo and Collins, 2010)

Why might we need higher-order dependency features? 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 $\langle \textit{scratch} \rightarrow \textit{with} \rangle$ and $\langle \textit{people} \rightarrow \textit{with} \rangle$, 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 grandchild 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. In general, higher-order features can yield substantial improvements in dependency parsing accuracy (e.g., Koo and Collins, 2010).

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. Projective third-order parsing can be performed in $\mathcal{O}(M^4)$ time and $\mathcal{O}(M^3)$ space (Koo and Collins, 2010). Approximate pruning algorithms can reduce this cost significantly by filtering out unpromising edges (Rush and Petrov, 2012).

Given the tractability of higher-order projective dependency parsing, you may be surprised to learn that non-projective second-order dependency parsing is NP-Hard! This can be proved by reduction from the vertex cover problem (Neuhaus and Bröker, 1997). One heuristic solution is to do projective parsing first, and then post-process the projective dependency parse to add non-projective edges (Nivre and Nilsson, 2005). More recent work has applied advanced techniques for approximate inference in graphical models, including belief propagation (Smith and Eisner, 2008), integer linear programming (Martins et al., 2009), variational inference (Martins et al., 2010), and Markov Chain Monte Carlo (Zhang et al., 2014).

10.3 Transition-based dependency parsing

Graph-based dependency parsing offers exact inference, meaning that it is possible to recover the best-scoring parse. But this exactness comes at a price: we can use only a limited set of features. These limitations are felt more keenly in dependency parsing than in sequence labeling, because second-order dependency features are critical to correctly identify certain types of attachments. Graph-based dependency parsing may also be criticized on the basis of intuitions about human language processing: people read and listen to sentences *sequentially*, incrementally building mental models of the sentence structure and meaning before getting to the end (Jurafsky, 1996). This seems hard to reconcile with graph-based algorithms, which perform bottom-up operations on the entire sentence, requiring the parser to keep every word in memory.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Transition-based algorithms address both of these objections. They work by moving through the sentence sequentially, while incrementally updating a stored representation of what has been read thus far. After processing the entire sentence, they return an analysis of its syntactic structure. A simple transition-based parser is the **arc-standard** parsing algorithm, which is similar to the LR algorithm that is used to parse programming languages. Transition-based parsing algorithms maintain a configuration state, which includes a stack where elements can be pushed and popped. They update the state incrementally through a series of actions, until the input is consumed and the stack is empty.

In the arc-standard parser, the configuration c is a tuple $c = (\sigma, \beta, A)$, where σ is a stack, β is the input buffer, and A is a set of dependency arcs. The initial state is $(\sigma = [0], \beta = w_{1:M}, A = \emptyset)$, where: $\sigma = [0]$ indicates that the stack begins with the root node; $\beta = w_{1:M}$ indicates that the buffer begins with the entire input string (indexed from 1); and $A = \emptyset$ means that there are not yet any arcs. We can then apply three possible actions:

- **SHIFT**: moves the first item from the input buffer on to the top of the stack,

$$(\sigma, i|\beta, A) \Rightarrow (\sigma|i, \beta, A), \quad [10.10]$$

where we write $i|\beta$ to indicate that i is the leftmost item in the input buffer, and $\sigma|i$ to indicate the result of pushing i on to stack σ .

- **ARC-LEFT**: creates a new left-facing arc between the item on the top of the stack and the first item in the input buffer. This item is then “popped” to the front of the input buffer, and the arc is added to A .

$$(\sigma|i, j|\beta, A) \Rightarrow (\sigma, j|\beta, A \cup (j, \ell, i)), \quad [10.11]$$

where ℓ is the (optional) label of the dependency arc.

- **ARC-RIGHT**: creates a new right-facing arc between the item on the top of the stack and the first item in the input buffer; this item is then “popped” to the front of the input buffer, and the arc is added to A .

$$(\sigma|i, j|\beta, A) \Rightarrow (\sigma, i|\beta, A \cup (i, \ell, j)), \quad [10.12]$$

where again ℓ is the label of the dependency arc.

The ARC-LEFT action cannot be performed when the root node 0 is on top of the stack, since this node must be the root of the entire tree. Neither ARC-LEFT nor ARC-RIGHT can be performed if the result would create a second incoming edge for any word. When the stack σ and the input buffer β are empty, parsing is complete.

(c) Jacob Eisenstein 2014-2017. Work in progress.

10.3.1 Learning transition-based parsers

Transition-based parsing requires selecting a series of actions. In parsing programming languages, shift-reduce parsers can choose the appropriate action deterministically, because programming languages are unambiguous by design. For natural language, we use machine learning classification to determine the best series of actions; for example, Yamada and Matsumoto (2003) use a support vector machine classifier (see § 1.2.1) to decide whether to shift or create a dependency arc at each stage in parsing.

To train a transition-based dependency parser, we can treat each parsing decision as a separate training instance. However, our ground truth input is not a list of parsing decisions, but rather, a dependency tree. We therefore require an **oracle** to convert the ground truth dependency tree into a list of parsing decisions, which can then be used as training data (Nivre, 2008).²

Typical features for transition-based dependency parsing include: the word and part-of-speech of the top element on the stack; the word and part-of-speech of the first, second, and third elements on the input buffer; pairs and triples of words and parts-of-speech from the top of the stack and the front of the buffer; the distance (in tokens) between the element on the top of the stack and the element in the front of the input buffer; the number of modifiers of each of these elements; and higher-order dependency features as described above in the section on graph-based dependency parsing. Zhang and Nivre (2011) describe a transition-based parser with rich features, which gave state-of-the-art performance (at the time) in both English and Chinese.

10.3.2 Pros and cons of transition-based dependency parsing

A key advantage of transition-based parsing is that it is much faster than graph-based methods. Since every word can be shifted once and every arc-creation action eliminates a word from the stack, the time complexity is linear in the length of the input. In contrast, graph-based parsing algorithms have quadratic or cubic time complexity.

Transition-based parsing can suffer from search errors, since an early mistake can make it impossible to find the best parse. This means that there could be an action sequence that would be preferred by the current parsing model, but is nonetheless not chosen because the first few actions in the sequence score badly. Put another way, transition-based parsing is **greedy** — unlike graph-based algorithms, which are guaranteed to find the best-scoring overall analysis. Solutions to this problem are discussed below.

Nonetheless, transition-based parsing achieves comparable accuracy to graph-based methods, in far less time (Nivre, 2004; Nivre et al., 2007). One reason is that in exchange

²**Spurious ambiguity** occurs when there are multiple derivations for the same dependency structure. This is the case in arc-standard dependency parsing: the structure $1 \leftarrow 2 \rightarrow 3$ can be obtained from two different action sequences (Cohen et al., 2012).

for giving up on global inference (and thereby accepting the possibility of search errors), we free ourselves from any restrictions on the features that can be used in the classifier that selects each parsing action. For example, features may consider any number of previous parsing decisions, any aspect of the current stack, and any part of the input.

10.3.3 Alternative transition-based parsing algorithms

Arc-eager dependency parsing changes the ARC-RIGHT action so that right dependents can be attached before all of **their** dependents have been found. In arc-eager parsing, the ARC-RIGHT action creates an arc, and then pushes both the parent and child elements on to the stack. To remove these elements, it adds an addition REDUCE action, which can be applied to elements on the stack for whom an incoming edge has already been identified. 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 A drawback of transition-based parsing is the possibility for search errors, in which a poor decision early in the parse will lead to cascading errors. **Beam search** is an improvement on greedy transition-based parsing, with the goal of eliminating search errors. As we move through the sentence, we keep a beam of possible hypotheses. For each element on the beam, we consider possible actions, and obtain a list of the top- k possibilities across all such actions. We then update the beam with the results of these top- k actions, and proceed (Zhang and Clark, 2008). Huang et al. (2012) offer alternative perceptron learning rules that are specifically designed for learning in the beam search setting.

Shift-reduce parsing for CFGs Transition-based parsing can also be used for parsing in (binarized) context-free grammars. Here we use a shift-reduce parser, where each reduce operation creates a new non-terminal that produces the top two elements in the stack. When the input is consumed and the only element on the stack is a tree derived from the start symbol S , the input has been completely parsed.

10.3.4 Neural transition-based parsing

[todo: 2-3 paragraphs about the following papers:]

- Each shift-reduce decision is made by a locally-trained neural network (Chen and Manning, 2014). See also (Dyer et al., 2015)
- Shift-reduce decisions are made by neural network trained on global conditional likelihood (Andor et al., 2016), using beam search.

(c) Jacob Eisenstein 2014-2017. Work in progress.



Figure 10.8: Google n-grams results for the bigram *write code* and the dependency arc *write => code* (and their morphological variants)

10.4 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*³

It is now possible to search Google n-grams by dependency edges; for example, finding the trend in how often a dependency edge has appeared over time. For example, we might be interested in knowing when people started talking about *writing code*, but we also want *write some code*, *write the code*, *write all the code*, etc. By searching on dependency edges, we can recover this information, as shown in Figure 10.8. This capability has implications for research in digital humanities, as shown by the analysis of Shakespeare performed by Muralidharan and Hearst (2013).

A classic application of dependency parsing is **relation extraction**, which is described in chapter 16. The goal of relation extraction is to identify entity pairs, such as

⟨TOLSTOY, WAR AND PEACE⟩
 ⟨MARQUÉZ, 100 YEARS OF SOLITUDE⟩
 ⟨SHAKESPEARE, A MIDSUMMER NIGHT’S DREAM⟩,

which stand in some relation to each other (in this case, the relation is authorship). Such entity pairs are often referenced via consistent chains of dependency relations. Therefore, dependency paths are often a useful feature in supervised systems which learn to detect

³Recall that the copula *is* is collapsed in many dependency grammars, such as the Universal Dependency treebank Nivre et al. (2016).

new instances of a relation, based on labeled examples of other instances of the same relation type (Culotta and Sorensen, 2004; Fundel et al., 2007; Mintz et al., 2009).

Cui et al. (2005) show how dependency parsing can improve question answering. For example, you might ask,

(10.1) *What % of the nation's cheese does Wisconsin produce?*

Now suppose your corpus contains this sentence:

(10.2) *In Wisconsin, where farmers produce 28% of the nation's cheese, ...*

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.

A final example comes from sentiment analysis. As discussed in chapter 3, the polarity of a sentence can be reversed by negation, e.g.

(10.3) *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, determining when key sentiment words are reversed (Wilson et al., 2005; Nakagawa et al., 2010).

Exercises

1. The dependency structure $1 \leftarrow 2 \rightarrow 3$, with 2 as the root, can be obtained from more than one set of actions in arc-standard parsing. List both sets of actions that can obtain this parse.

Part III

Meaning

Chapter 11

Compositional semantics

A grand ambition of natural language processing is to convert natural language into a representation that supports inferences about meaning. Indeed, many potential applications of language technology involve some level of semantic understanding:

- Answering questions, such as *where is the nearest coffeeshop?* or *what is the middle name of the mother of the 44th President of the United States?*.
- Building a robot that can follow natural language instructions to execute tasks.
- Translating a sentence from one language into another, while preserving the underlying meaning.
- Fact-checking an article by searching the web for contradictory evidence.
- Logic-checking an argument by identifying contradictions, ambiguity, and unsupported assertions.

Each of these applications can be performed by converting natural language into a **meaning representation**. To be useful, a meaning representation must meet several criteria:

- **c1**: it should be unambiguous (unlike natural language);
- **c2**: it should provide a way to link language to external knowledge, observations, and actions;
- **c3**: it should support computational **inference**.

Much more than this can be said about the question of how best to represent knowledge for computation (e.g., Sowa, 2000), but we will focus on these three criteria.

11.1 Meaning and denotation

The first criterion for a meaning representation is that statements in the representation should be unambiguous — they should have only one possible interpretation. Natural language does not have this property: as we saw in chapter 9, sentences like *cats scratch people with claws* have multiple interpretations.

But what does it mean for a statement to be unambiguous? Programming languages provide a useful example: the output of a program is completely specified by the rules of the language, and the properties of the environment in which the program is run. For example, the python code `5 + 3` will have the output 8, as will the codes `(4*4) - (3*3) + 1` and `((8))`. This output is known as the **denotation** of the program, and can be written as,

$$\llbracket 5+3 \rrbracket = \llbracket (4*4) - (3*3) + 1 \rrbracket = \llbracket ((8)) \rrbracket = 8. \quad [11.1]$$

In this sense, the program’s output is its “meaning”.

The denotations of these arithmetic expressions are determined by the meaning of the **constants** (e.g., 5, 3) and the **relations** (e.g., +, *, (,)). Now let’s consider another snippet of python code, `double(4)`. The denotation of this code could be, $\llbracket \text{double}(4) \rrbracket = 8$, or it could be $\llbracket \text{double}(4) \rrbracket = 44$ — it depends on the meaning of `double`. This meaning is defined in a **world model** \mathcal{M} as an infinite set of pairs. We write the denotation with respect to model \mathcal{M} as $\llbracket \cdot \rrbracket_{\mathcal{M}}$, e.g., $\llbracket \text{double} \rrbracket_{\mathcal{M}} = \{\langle 0, 0 \rangle, \langle 1, 2 \rangle, \langle 2, 4 \rangle, \dots\}$. The world model would also define the (infinite) list of constants, e.g., $\{0, 1, 2, \dots\}$. As long as the denotation of string ϕ in model \mathcal{M} can be computed unambiguously, the language can be said to be unambiguous.

This approach to meaning is known as **model-theoretic semantics**, and it addresses not only criterion *c1* (no ambiguity), but also *c2* (connecting language to external knowledge, observations, and actions). For example, we can connect a representation of the meaning of a statement like *the capital of Georgia* with a world model that includes knowledge base of geographical facts, obtaining the denotation *Atlanta*. We might populate a world model by applying an image analysis algorithm to Figure 11.1, and then use this world model to evaluate **propositions** like *a man is riding a moose*. Another desirable property of model-theoretic semantics is that when the facts change, the denotations change too: the meaning representation of *President of the USA* would have a different denotation in the model \mathcal{M}_{2014} as it would in \mathcal{M}_{2022} .

11.2 Logical representations of meaning

If we can find a meaning representation which supports model-theoretic denotation, then we will have met criteria *c1* and *c2*. The final criterion is *c3*, which requires that the meaning representation support inference — for example, automatically deducing new



Figure 11.1: A (doctored) image, which could be the basis for a world model [todo: the image is from 1912, so out of copyright? <https://blogs.harvard.edu/houghton/2013/09/20/myths-debunked-sadly-theodore-roosevelt-never-rode-a-moose/>]

facts from known premises. While many representations have been proposed that meet these criteria, the most mature is the language of first-order logic.¹

11.2.1 Propositional logic

The bare bones of logical meaning representation are Boolean operations on propositions:

Propositional symbols Greek symbols like ϕ and ψ will be used to represent **propositions**, which are statements that are either true or false. For example, ϕ may correspond to the proposition, *bagels are delicious*.

Boolean operators We can build up more complex propositional formulas from Boolean operators. These include:

- Negation $\neg\phi$, which is true if ϕ is false.
- Conjunction, $\phi \wedge \psi$, which is true if both ϕ and ψ are true.

¹Relevant alternatives include the “variable-free” representation used in semantic parsing of geographical queries (Zelle and Mooney, 1996) and robotic control (Ge and Mooney, 2005), and dependency-based compositional semantics (Liang et al., 2013).

- Disjunction, $\phi \vee \psi$, which is true if at least one of P and Q is true
- Implication, $\phi \Rightarrow \psi$, which is true unless ϕ is true and ψ is false. Implication has identical truth conditions to $\neg\phi \vee \psi$.
- Equivalence, $\phi \Leftrightarrow \psi$, which is true if ϕ and ψ are both true or both false. Equivalence has identical truth conditions to $(\phi \Rightarrow \psi) \wedge (\psi \Rightarrow \phi)$.

It is not strictly necessary to have all five Boolean operators: readers familiar with Boolean logic will know that it is possible to construct all other operators from either the NAND (not-and) or NOR (not-or) operators. Nonetheless, it is typical to use all five operators above for clarity. It is possible to define a number of “laws” for these Boolean operators, such as,

- **commutativity:** $\phi \wedge \psi = \psi \wedge \phi$, $\phi \vee \psi = \psi \vee \phi$
- **associativity:** $\phi \wedge (\psi \wedge \chi) = (\phi \wedge \psi) \wedge \chi$, $\phi \vee (\psi \vee \chi) = (\phi \vee \psi) \vee \chi$
- **complementation:** $\phi \wedge \neg\phi = \perp$, $\phi \vee \neg\phi = \top$, where \top indicates a true proposition and \perp indicates a false proposition.

These laws can be combined to derive further equivalences, which can support logical inferences. For example, suppose $\phi = \text{The music is loud}$ and $\psi = \text{Max can't sleep}$. Then if we are given,

$$\begin{array}{ll} \phi \Rightarrow \psi & \text{If the music is loud, Max can't sleep.} \\ \phi & \text{The music is loud.} \end{array}$$

we can derive ψ (*Max can't sleep*) by application of **modus ponens**, which is one of a set of **inference rules** that can be derived from more basic laws and used to manipulate propositional formulas. **Automated theorem provers** are capable of applying inference rules to a set of premises to derive desired propositions (Loveland, 2016).

11.2.2 First-order logic

Propositional logic is so named because it treats propositions as its base units. However, we would also like to reason about the content of the propositions themselves. For example,

(11.1) If anyone is making noise, then Max can't sleep.

(11.2) Abigail is making noise.

To understand the relationship between the statement *anyone is making noise* and the statement *Abigail is making noise*, we need some additional formal machinery. This is provided by **first-order logic** (FOL).

(c) Jacob Eisenstein 2014-2017. Work in progress.

In FOL, logical propositions can be constructed from relationships between entities. Specifically, FOL extends propositional logic with the following classes of terms:

Constants These are elements that name individual entities in the model, such as MAX and ABIGAIL. The **denotation** of each constant in a model \mathcal{M} is an element in the model, e.g., $\llbracket \text{MAX} \rrbracket = m$ and $\llbracket \text{ABIGAIL} \rrbracket = a$.

Relations Relations can be thought of as sets of entities, or sets of tuples. For example, the relation CAN-SLEEP is defined as the set of entities who can sleep, and has the denotation $\llbracket \text{CAN-SLEEP} \rrbracket = \{a, m, \dots\}$. We can test the truth value of the proposition CAN-SLEEP(MAX) by asking whether $\llbracket \text{MAX} \rrbracket \in \llbracket \text{CAN-SLEEP} \rrbracket$. Logical relations that are defined over sets of entities are sometimes called **properties**.

Relations may also be ordered tuples of entities. For example BROTHER(MAX,ABIGAIL) expresses the proposition that MAX is the brother of ABIGAIL. The denotation of such relations is a set of tuples, $\llbracket \text{BROTHER} \rrbracket = \{(m, a), (x, y), \dots\}$. We can test the truth value of the proposition BROTHER(MAX,ABIGAIL) by asking whether the tuple $(\llbracket \text{MAX} \rrbracket, \llbracket \text{ABIGAIL} \rrbracket)$ is in the denotation $\llbracket \text{BROTHER} \rrbracket$.

We can now express statements like *Max can't sleep* and *Max is Abigail's brother*:

$$\neg \text{CAN-SLEEP}(\text{MAX})$$

$$\text{BROTHER}(\text{MAX}, \text{ABIGAIL}).$$

We can also combine these statements using Boolean operators, such as,

$$(\text{BROTHER}(\text{MAX}, \text{ABIGAIL}) \vee \text{BROTHER}(\text{MAX}, \text{STEVE})) \Rightarrow \neg \text{CAN-SLEEP}(\text{MAX}).$$

We have thus far described a fragment of first-order logic, which permits only statements about specific entities. To support inferences about statements like *If anyone is making noise, then Max can't sleep*, we will need two additional elements in the meaning representation:

Variables These are mechanisms for referring to entities that are not locally specified. We can then write CAN-SLEEP(x) or BROTHER(x , ABIGAIL). In these cases, x is an **free variable**, meaning that we have not committed to any particular assignment.

Quantifiers Variables are bound by quantifiers. There are two quantifiers in first-order logic.²

²In first-order logic, it is possible to quantify only over entities. In **second-order logic**, it is possible to quantify over properties, supporting statements like *Butch has every property that a good boxer has*,

$$\forall P \forall x ((\text{GOOD-BOXER}(x) \Rightarrow P(x)) \Rightarrow P(\text{BUTCH})) \quad [11.2]$$

This example is from Blackburn and Bos (2005), who also show how first-order logic can “approximate” second-order and higher-order logics, by reifying sets as additional entities in the model.

- The **existential quantifier** \exists , which indicates that there must be at least one entity to which the variable can refer. For example, the statement $\exists x \text{MAKES-NOISE}(x)$ indicates that there is at least one entity for which **MAKES-NOISE** is true.
- The **universal quantifier** \forall , which indicates that the variable must be able to refer to any entity. For example, the statement,

$$\neg \text{MAKES-NOISE}(\text{ABIGAIL}) \Rightarrow (\forall x \text{CAN-SLEEP}(x)) \quad [11.3]$$

asserts that every entity can sleep if Abigail does not make noise.

The expressions $\exists x$ and $\forall x$ make x into a **bound variable**. A formula that contains no free variables is a **sentence**.

Functions Functions map from entities to entities, e.g., $\llbracket \text{CAPITAL-OF}(\text{GEORGIA}) \rrbracket = \llbracket \text{ATLANTA} \rrbracket$. With functions, we also add an equality operator, so that it is possible to make statements like,

$$\forall x \exists y \text{MOTHER-OF}(x) = \text{DAUGHTER-OF}(y). \quad [11.4]$$

Note that **MOTHER-OF** is a functional analogue of the relation **MOTHER**, so that $\text{MOTHER-OF}(x) = y$ if $\text{MOTHER}(x, y)$. Any logical formula that uses functions can be rewritten using only relations and quantification. For example,

$$\text{MAKES-NOISE}(\text{MOTHER-OF}(\text{ABIGAIL})) \quad [11.5]$$

can be rewritten as $\exists x \text{MAKES-NOISE}(x) \wedge \text{MOTHER}(x, \text{ABIGAIL})$.

An important property of quantifiers is that the order can matter. Unfortunately, natural language is rarely clear about this! The issue is demonstrated by examples like *everyone speaks a language*, which has the following interpretations:

$$\forall x \exists y \text{SPEAKS}(x, y) \quad [11.6]$$

$$\exists y \forall x \text{SPEAKS}(x, y). \quad [11.7]$$

In the first case, y may refer to several different languages, while in the second case, there is a single y that is spoken by everyone.

Truth-conditional semantics

One way to look at the meaning of an FOL sentence ϕ is as a set of **truth conditions**, or models under which ϕ is satisfied. But how can we determine whether a sentence in first-order logic is true or false? We will approach this inductively, starting with a predicate applied to a tuple of constants. The truth of such a sentence depends on whether the

(c) Jacob Eisenstein 2014-2017. Work in progress.

tuple of denotations of the constants is in the denotation of the predicate. For example, $\text{CAPITAL}(\text{GEORGIA}, \text{ATLANTA})$ is true in model \mathcal{M} iff,

$$\langle \llbracket \text{GEORGIA} \rrbracket_{\mathcal{M}}, \llbracket \text{ATLANTA} \rrbracket_{\mathcal{M}} \rangle \in \llbracket \text{CAPITAL} \rrbracket_{\mathcal{M}}. \quad [11.8]$$

The Boolean operators \wedge, \vee, \dots provide ways to construct more complicated sentences, and the truth of such statements can be assessed based on the truth tables associated with these operators. The statement $\exists x\phi$ is true if there is some assignment of the variable x to an entity in the model such that ϕ is true; the statement $\forall x\phi$ is true if ϕ is true under all possible assignments of x . More formally, we would say that ϕ is **satisfied** under \mathcal{M} , written as $\mathcal{M} \models \phi$.

Truth conditional semantics allows us to define several other properties of sentences and pairs of sentences. Suppose that in every \mathcal{M} under which ϕ is satisfied, another formula ψ is also satisfied; then ϕ **entails** ψ , sometimes written $\phi \models \psi$ [**todo: double check**]. For example,

$$\text{CAPITAL}(\text{GEORGIA}, \text{ATLANTA}) \models \exists x \text{CAPITAL}(\text{GEORGIA}, x). \quad [11.9]$$

A statement that is satisfied under any model, such as $\phi \vee \neg\phi$, is **valid**; a statement that is not satisfied under any model, such as $\phi \wedge \neg\phi$, is **unsatisfiable**, or **inconsistent**. A **model checker** is a program that determines whether a sentence ϕ is satisfied in \mathcal{M} . A **model builder** is a program that constructs a model in which ϕ is satisfied. The problems of checking for consistency and validity in first-order logic are **undecidable**, meaning that there is no algorithm that can automatically determine whether an FOL formula is valid or inconsistent.

Inference in first-order logic

Our original goal was to support inferences that combine general statements *If anyone is making noise, then Max can't sleep* with specific statements like *Abigail is making noise*. We can now represent such statements in first-order logic, but how are we to perform the inference that *Max can't sleep*? One approach is to use “generalized” versions of propositional inference rules like modus ponens, which can be applied to FOL formulas. By repeatedly applying such inference rules to a knowledge base of facts, it is possible to produce proofs of desired propositions. To find the right sequence of inferences to derive a desired theorem, classical artificial intelligence search algorithms like backward chaining can be applied. Such algorithms are implemented in interpreters for the `prolog` logic programming language (Pereira and Shieber, 2002).

11.3 Semantic parsing and the lambda calculus

The previous section has laid out a lot of formal machinery, which we will now try to unite with natural language. Given an English sentence like *Alex likes Brit*, how can we

(c) Jacob Eisenstein 2014-2017. Work in progress.

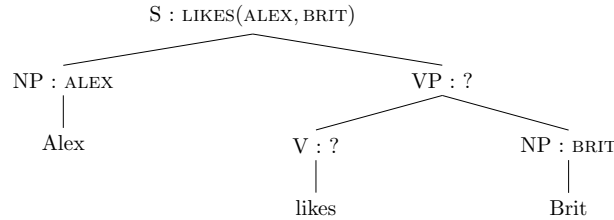


Figure 11.2: The principle of compositionality requires that we identify meanings for the constituents *likes* and *likes Brit* that will make it possible to compute the meaning for the entire sentence.

obtain the desired first-order logical representation, $\text{LIKES}(\text{ALEX}, \text{BRIT})$? This is the task of **semantic parsing**. Just as a syntactic parser is a function from a natural language sentence to a syntactic structure such as a phrase structure tree, a semantic parser is a function from natural language to logical formulas.

As in syntactic analysis, semantic parsing is difficult because the space of inputs and outputs is very large, and their interaction is complex. Our best hope is that, like syntactic parsing, semantic parsing can somehow be decomposed into simpler sub-problems. This idea, usually attributed to the German philosopher Gottlob Frege, is called the **principle of compositionality**: the meaning of a complex expression is a function of the meanings of that expression’s constituent parts. We will define these “constituent parts” as syntactic elements like noun phrases and verb phrases. These constituents are combined using function application: if the syntactic parse contains the production $x \rightarrow y z$, then the semantics of x , written $x.\text{sem}$, will be computed as a function of the semantics of the constituents, $y.\text{sem}$ and $z.\text{sem}$.^{3 4}

11.3.1 The lambda calculus

Let’s see how this works for a simple sentence like *Alex likes Brit*, whose syntactic structure is shown in Figure 11.2. Our goal is the formula, $\text{LIKES}(\text{ALEX}, \text{BRIT})$, and it is clear that the meaning of the constituents *Alex* and *Brit* should be ALEX and BRIT . That leaves two more constituents: the verb *likes*, and the verb phrase *likes Brit*. How can we define the meaning of these units in a way that enables us to recover the desired meaning for the

³§ 8.3.1 briefly discusses alternative syntactic formalisms, including Combinatory Categorical Grammar (CCG). CCG is argued to be particularly well-suited to semantic parsing (Hockenmaier and Steedman, 2007), and is used in much of the contemporary work on machine learning for semantic parsing, summarized in § 11.4.

⁴The approach of algorithmically building up meaning representations from a series of operations on the syntactic structure of a sentence is generally attributed to the philosopher Richard Montague, who published a series of influential papers on the topic in the early 1970s (e.g., Montague, 1973).

entire sentence? If the meanings of *Alex* and *Brit* are constants, then the meanings of *likes* and *likes Brit* must be functional expressions, which can be applied to their siblings to produce the desired analyses.

Modeling these partial analyses requires extending the first-order logic meaning representation. We do this by adding **lambda expressions**, which are descriptions of anonymous functions,⁵ e.g.,

$$\lambda x. \text{LIKES}(x, \text{BRIT}). \quad [11.10]$$

We take this functional expression to be the meaning of the verb phrase *likes Brit*; it takes a single argument, and returns the result of substituting that argument for x in the expression $\text{LIKES}(x, \text{BRIT})$. We write this substitution as,

$$(\lambda x. \text{LIKES}(x, \text{BRIT}))@ \text{ALEX} = \text{LIKES}(\text{ALEX}, \text{BRIT}), \quad [11.11]$$

with the symbol “@” indicating function application. Function application in the lambda calculus is sometimes called **β -reduction** or **β -conversion**. We will write $\phi@ \psi$ to indicate a function application to be performed by β -reduction, and $\phi(\psi)$ to indicate a function or predicate in the final logical form.

Equation 11.11 shows how to obtain the desired semantics for the sentence *Alex likes Brit*: by applying the lambda expression $\lambda x. \text{LIKES}(x, \text{BRIT})$ to the logical constant ALEX . This rule of composition can be specified in a syntactic-semantic grammar: for the syntactic production $S \rightarrow \text{NP VP}$, we have the semantic rule $\text{VP.sem}@ \text{NP.sem}$.

The meaning of the meaning of the transitive verb phrase can also be obtained by function application on its syntactic constituents. For the syntactic production $\text{VP} \rightarrow \text{V NP}$, we apply the semantic rule,

$$\text{VP.sem} = (\text{V.sem})@ \text{NP.sem} \quad [11.12]$$

$$= (\lambda y. \lambda x. \text{LIKES}(x, y))@ (\text{BRIT}) \quad [11.13]$$

$$= \lambda x. \text{LIKES}(x, \text{BRIT}). \quad [11.14]$$

Here we have defined the meaning of the transitive verb *likes* as a lambda expression whose output is **another** lambda expression: it takes y as an argument to fill in one of the slots in the LIKES relation, and returns a lambda expression that is ready to take an argument to fill in the other slot.⁶

⁵Formally, all first-order logic formulas are lambda expressions; in addition, if ϕ is a lambda expression, then $\lambda x. \phi$ is also a lambda expression. Readers who are familiar with functional programming will recognize lambda expressions from their use in programming languages such as Lisp and Python.

⁶This can be written in a few different ways. More informally, we can write $\lambda y. x. \text{LIKES}(x, y)$, indicating a lambda expression that takes two arguments; this would be acceptable in functional programming. More formally, logicians (e.g., Carpenter, 1997) often write $\lambda y. \lambda x. \text{LIKES}(x)(y)$, indicating that each lambda expression takes exactly one argument.

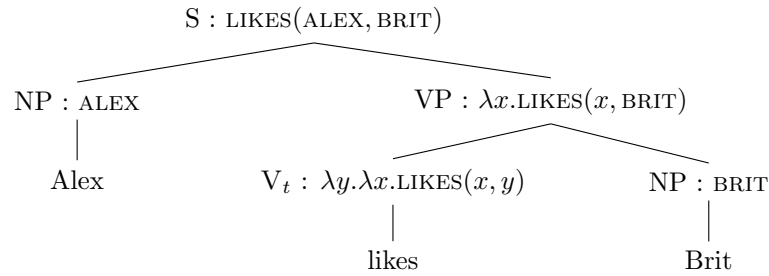


Figure 11.3: Derivation of the semantic representation for *Alex likes Brit* in the grammar G_1 .

S	→ NP VP	VP.sem@NP.sem
VP	→ V_t NP	V_t .sem@NP.sem
VP	→ V_i	V_i .sem
V_t	→ <i>likes</i>	$\lambda y.\lambda x.LIKES(x, y)$
V_i	→ <i>sleeps</i>	$\lambda x.SLEEPS(x)$
NP	→ <i>Alex</i>	ALEX
NP	→ <i>Brit</i>	BRIT

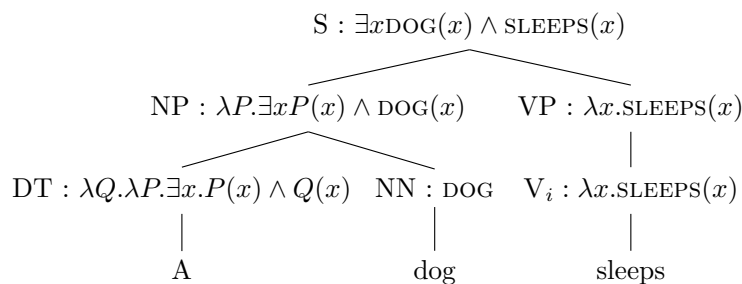
Table 11.1: G_1 , a minimal syntactic/semantic context-free grammar

Table 11.1 shows a minimal syntactic/semantic grammar fragment, which we will call G_1 . The complete **derivation** of *Alex likes Brit* in G_1 is shown in Figure 11.3. In addition to the transitive verb *likes*, the grammar also includes the intransitive verb *sleeps*; it should be clear how to derive the meaning of sentences like *Alex sleeps*. For verbs that can be either transitive or intransitive, such as *eats*, we would have two terminal productions, one for each sense (terminal productions are also called the **lexical entries**). Indeed, most of the grammar is in the **lexicon** (the terminal productions), since these productions select the basic units of the semantic interpretation.

11.3.2 Quantification

Things get more complicated when we move from sentences about named entities to sentences that involve more general noun phrases. Let's consider the example, *A dog sleeps*, which has the meaning $\exists x \text{DOG}(x) \wedge \text{SLEEPS}(x)$. Clearly, the DOG relation will be introduced by the word *dog*, and the SLEEP relation will be introduced by the word *sleeps*. The existential quantifier \exists must be introduced by the lexical entry for the determiner *a*.⁷

⁷Conversely, the sentence *Every dog sleeps* would involve a universal quantifier, $\forall x \text{DOG}(x) \Rightarrow \text{SLEEPS}(x)$. The definite article *the* requires more consideration, since *the dog* must refer to some dog which is uniquely

Figure 11.4: Derivation of the semantic representation for *A dog sleeps*, in grammar G_2

However, this seems problematic for the compositional approach taken in the grammar G_1 : if the semantics of the noun phrase *a dog* is an existentially quantified expression, how can it be the argument to the semantics of the verb *sleeps*, which expects an entity? And where does the logical conjunction come from?

There are a few different approaches to handling these issues.⁸ We will begin by reversing the semantic relationship between subject NPs and VPs, so that the production $S \rightarrow \text{NP VP}$ has the semantics $\text{NP.sem} @ \text{VP.sem}$: the meaning of the sentence is now the semantics of the noun phrase applied to the verb phrase. The implications of this change are best illustrated by exploring the derivation of the example, shown in Figure 11.4. Let's start with the indefinite article *a*, to which we assign the rather intimidating semantics,

$$\lambda P. \lambda Q. \exists x P(x) \wedge Q(x). \quad [11.15]$$

This is a lambda expression that takes two **relations** as arguments, P and Q . The relation P is scoped to the outer lambda expression, so it will be provided by the immediately adjacent noun, which in this case is **DOG**. Thus, the noun phrase *a dog* has the semantics,

$$\text{NP.sem} = \text{DET.sem} @ \text{NN.sem} \quad [11.16]$$

$$= (\lambda P. \lambda Q. \exists x P(x) \wedge Q(x)) @ (\text{DOG}) \quad [11.17]$$

$$= \lambda Q. \exists x \text{DOG}(x) \wedge Q(x). \quad [11.18]$$

This is a lambda expression that is expecting another relation, Q , which will be provided by the verb phrase, **SLEEPS**. This gives the desired analysis, $\exists x \text{DOG}(x) \wedge \text{SLEEPS}(x)$.⁹

identifiable, perhaps from contextual information external to the sentence. Carpenter (1997, pp. 96-100) summarizes recent approaches to handling definite descriptions.

⁸Carpenter (1997) offers an alternative treatment based on combinatory categorial grammar.

⁹When applying β -reduction to arguments that are themselves lambda expressions, be sure to use unique variable names to avoid confusion. For example, it is important to distinguish the x in the semantics for *a* from the x in the semantics for *likes*. Variable names are abstractions, and can always be changed — this is known as α -conversion. For example, $\lambda x. P(x)$ can be converted to $\lambda y. P(y)$, etc.

If noun phrases like *a dog* are interpreted as lambda expressions, then proper nouns like *Alex* must be treated in the same way. This is achieved by **type-raising** from constants to lambda expressions, $x \Rightarrow \lambda P.P(x)$. After type-raising, the semantics of *Alex* is $\lambda P.P(\text{ALEX})$ — a lambda expression that expects a relation to tell us something about ALEX.¹⁰ Make sure you see how the analysis in Figure 11.4 can be applied to the sentence *Alex sleeps*.

To handle direct objects, we will perform the same type-raising operation on transitive verbs: the meaning of verbs such as *likes* will be raised to,

$$\lambda P.\lambda x.P(\lambda y.\text{LIKES}(x, y)) \quad [11.19]$$

As a result, we can keep the verb phrase production $\text{VP.sem} = \text{V.sem} @ \text{NP.sem}$, knowing that the direct object will provide the function P in Equation 11.19. To see how this works, let's analyze the verb phrase *likes a dog*. After uniquely relabeling each lambda variable, we have,

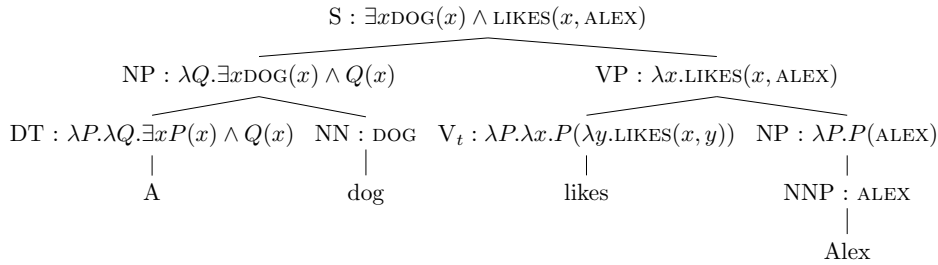
$$\begin{aligned} \text{VP.sem} &= \text{V.sem} @ \text{NP.sem} \\ &= (\lambda P.\lambda x.P(\lambda y.\text{LIKES}(x, y))) @ (\lambda Q.\exists z \text{DOG}(z) \wedge Q(z)) \\ &= \lambda x.(\lambda Q.\exists z \text{DOG}(z) \wedge Q(z)) @ (\lambda y.\text{LIKES}(x, y)) \\ &= \lambda x.\exists z \text{DOG}(z) \wedge (\lambda y.\text{LIKES}(x, y)) @ z \\ &= \lambda x.\exists z \text{DOG}(z) \wedge \text{LIKES}(x, z). \end{aligned}$$

These changes are summarized in the revised grammar G_2 , shown in Table 11.2. Figure 11.5 shows a derivation that involves a transitive verb, an indefinite noun phrase, and a proper noun.

11.4 Learning semantic parsers

As with syntactic parsing, any syntactic/semantic grammar with sufficient coverage will **overgenerate**, producing many possible analyses for any given sentence. Machine learning is the dominant approach to selecting a single analysis. We will focus on algorithms that learn to score logical forms by attaching weights to features of their derivations (Zettlemoyer and Collins, 2005). Alternative approaches include transition-based parsing (Zelle

¹⁰Compositional semantic analysis is often supported by **type systems**, which make it possible to check whether a given function application is valid. The base types are entities e and truth values t . A property, such as DOG , is a function from entities to truth values, so its type is written $\langle e, t \rangle$. A transitive verb has type $\langle e, \langle e, t \rangle \rangle$: after receiving the first entity (the direct object), it returns a function from entities to truth values, which will be applied to the subject of the sentence. The type-raising operation $x \Rightarrow \lambda P.P(x)$ corresponds to a change in type from e to $\langle \langle e, t \rangle, t \rangle$: it expects a function from entities to truth values, and returns a truth value.

Figure 11.5: Derivation of the semantic representation for *A dog likes Alex*.

S	→ NP VP	NP.sem@VP.sem
VP	→ V _t NP	V _t .sem@NP.sem
VP	→ V _i	V _i .sem
NP	→ DET NN	DET.sem@NN.sem
NP	→ NNP	λP.P(NNP.sem)
DET	→ <i>a</i>	λP.λQ.∃x P(x) ∧ Q(x)
DET	→ <i>every</i>	λP.λQ.∀x (P(x) ⇒ Q(x))
V _t	→ <i>likes</i>	λP.λx.P(λy.LIKES(x, y))
V _i	→ <i>sleeps</i>	λx.SLEEPS(x)
NN	→ <i>dog</i>	DOG
NNP	→ <i>Alex</i>	ALEX
NNP	→ <i>Brit</i>	BRIT

Table 11.2: G_2 , a syntactic/semantic context-free grammar fragment, which supports quantified noun phrases

and Mooney, 1996; Misra and Artzi, 2016) and methods inspired by machine translation (Wong and Mooney, 2006). Methods also differ in the form of supervision used for learning, which can range from complete derivations to much more limited training signals. We will begin with the case of complete supervision, and then consider how learning is still possible even when seemingly key information is missing.

Datasets Early work on semantic parsing focused on geographical queries, such as *What states border Texas*. The GeoQuery dataset of Zelle and Mooney (1996) was originally coded in prolog, but has subsequently been expanded and converted into the SQL database query language by Popescu et al. (2003) and into first-order logic with lambda calculus by Zettlemoyer and Collins (2005), providing logical forms like $\lambda x. \text{STATE}(x) \wedge \text{BORDERS}(x, \text{TEXAS})$. Another early dataset consists of instructions for RoboCup robot soccer teams (Kate et al., 2005). More recent work has focused on broader domains, such as the Freebase database (Bol-

lacker et al., 2008), for which queries have been annotated by Krishnamurthy and Mitchell (2012) and Cai and Yates (2013), as well on child-directed speech (Kwiatkowski et al., 2012) and elementary school science exams (Krishnamurthy, 2016).

11.4.1 Learning from derivations

Let $w^{(i)}$ indicate a sequence of text, and let $y^{(i)}$ indicate the desired logical form. For example:

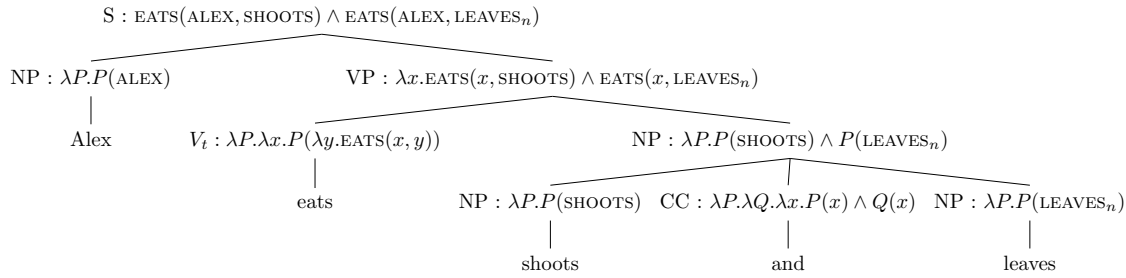
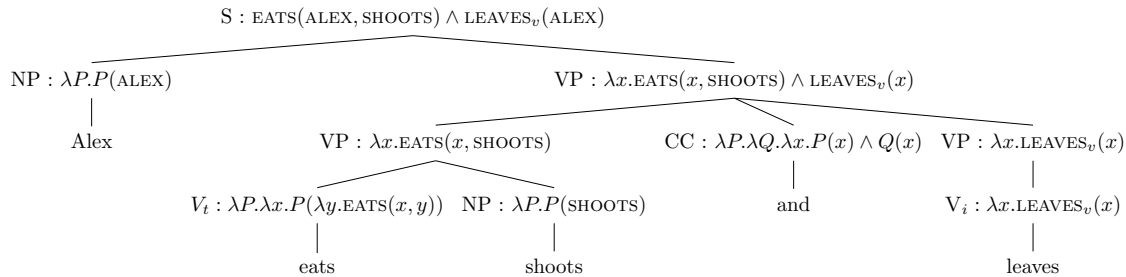
$$\begin{aligned} w^{(i)} &= \text{Alex eats shoots and leaves} \\ y^{(i)} &= \text{EATS(ALEX,SHOOTS)} \wedge \text{EATS(ALEX,LEAVES)} \end{aligned}$$

In the standard supervised learning paradigm that was introduced in § 1.2, we first define a feature function, $f(w, y)$, and then learn weights on these features, so that $y^{(i)} = \operatorname{argmax}_y \theta \cdot f(w, y)$. The weight vector θ is learned by comparing the features of the true label $f(w^{(i)}, y^{(i)})$ against either the features of the predicted label $f(w^{(i)}, \hat{y})$ (perceptron, support vector machine) or the expected feature vector $E_{y|w}[f(w^{(i)}, y)]$ (logistic regression).

While this basic framework seems similar to discriminative syntactic parsing, there is a crucial difference. In (context-free) syntactic parsing, the annotation $y^{(i)}$ contains all of the syntactic productions; indeed, the task of identifying the correct set of productions is identical to the task of identifying the syntactic structure. In semantic parsing, this is not the case: the logical form $\text{EATS(ALEX,SHOOTS)} \wedge \text{EATS(ALEX,LEAVES)}$ does not reveal the syntactic/semantic productions that were used to obtain it. Indeed, there may be **spurious ambiguity**, so that a single logical form can be reached by multiple derivations. (We previously encountered spurious ambiguity in transition-based dependency parsing, § 10.3.1.)

Let us introduce an additional variable z , representing the **derivation** of the logical form y from the text w . We assume that the feature function decomposes across the productions in the derivation, $f(w, z, y) = \sum_{t=1}^T f(w, z_t, y)$, where z_t indicates a single syntactic/semantic production. For example, we might have a feature for the high-level production $S \rightarrow \text{NP VP} : \text{NP.sem@VP.sem}$, as well as for terminal productions like $\text{NNP} \rightarrow \text{Alex} : \text{ALEX}$. Under this decomposition, we can compute scores for each semantically-annotated subtree in the analysis of w , and can therefore apply bottom-up parsing algorithms like CKY (§ 9.1) to find the best-scoring semantic analysis.

Figure 11.6 shows a derivation of the correct semantic analysis of the sentence *Alex eats shoots and leaves*, in a simplified grammar in which the plural noun phrases *shoots* and *leaves* are interpreted as logical constants SHOOTS and LEAVES_n . Figure 11.7 shows a derivation of an incorrect analysis. Assuming one feature per production, the perceptron update is shown in Table 11.3. From this update, the parser would learn to prefer the noun

Figure 11.6: Derivation for gold semantic analysis of *Alex eats shoots and leaves*Figure 11.7: Derivation for incorrect semantic analysis of *Alex eats shoots and leaves*

$NP_1 \rightarrow NP_2 \text{ CC } NP_3$	$(CC.sem@(NP_2.sem))@(NP_3.sem)$	+1
$VP_1 \rightarrow VP_2 \text{ CC } VP_3$	$(CC.sem@(VP_2.sem))@(VP_3.sem)$	-1
$NP \rightarrow leaves$	$LEAVES_n$	+1
$VP \rightarrow V_i$	$V_i.sem$	-1
$V_i \rightarrow leaves$	$\lambda x.LEAVES_v$	-1

Table 11.3: Perceptron update for analysis in Figure 11.6 (gold) and Figure 11.7 (predicted)

interpretation of *leaves* over the verb interpretation. It would also learn to prefer noun phrase coordination over verb phrase coordination. Note that we could easily replace the perceptron with a conditional random field. In this case, the online updates would be based on feature expectations, which can be computed using bottom-up algorithms like inside-outside (§ 9.5).

(c) Jacob Eisenstein 2014-2017. Work in progress.

11.4.2 Learning from logical forms

Complete derivations are expensive to annotate, and are rarely available.¹¹ More recent work has focused on learning from logical forms directly, while treating the derivations as **latent variables** (Zettlemoyer and Collins, 2005). In a conditional probabilistic model over logical forms y and derivations z , we have,

$$p(y, z | w) = \frac{\exp(\theta \cdot f(w, z, y))}{\sum_{y', z'} \exp(\theta \cdot f(w, z', y'))}, \quad [11.20]$$

which is the standard log-linear model, applied to the logical form y and the derivation z .

Since the derivation z completely determines the logical form y , it may seem silly to model the joint probability over y and z . However, since z is unknown, it can be marginalized out,

$$p(y | w) = \sum_z p(y, z | w). \quad [11.21]$$

We can then have the semantic parser select the logical form with the maximum log marginal probability,

$$\log \sum_z p(y, z | w) = \log \sum_z \frac{\exp(\theta \cdot f(w, z, y))}{\sum_{y', z'} \exp(\theta \cdot f(w, z', y'))} \quad [11.22]$$

$$\propto \log \sum_z \exp(\theta \cdot f(w, z, y)) \quad [11.23]$$

$$\geq \max_z \theta \cdot f(w, z, y). \quad [11.24]$$

Note that it is impossible to push the log term inside the sum over z , meaning that our usual linear scoring function does not apply. We can recover this scoring function only in approximation, by taking the max (rather than the sum) over derivations z , which provides a lower bound.

Learning can be performed by maximizing the log marginal likelihood,

$$\ell(\theta) = \sum_{i=1}^N \log p(y^{(i)} | w^{(i)}; \theta) \quad [11.25]$$

$$= \sum_{i=1}^N \log \sum_z p(y^{(i)}, z^{(i)} | w^{(i)}; \theta). \quad [11.26]$$

¹¹An exception is the work of Ge and Mooney (2005), who annotate the meaning of each syntactic constituents for several hundred sentences.

This log-likelihood is not **convex** in θ , unlike the log-likelihood of a fully-observed conditional random field. This means that learning can give different results depending on the initialization.

The derivative of Equation 11.26 is,

$$\frac{\partial \ell_i}{\partial \theta} = \sum_z p(z \mid \mathbf{y}, \mathbf{w}; \theta) \mathbf{f}(\mathbf{w}, z, \mathbf{y}) - \sum_{\mathbf{y}', z'} p(\mathbf{y}', z' \mid \mathbf{w}; \theta) \mathbf{f}(\mathbf{w}, z', \mathbf{y}') \quad [11.27]$$

$$= E_{z \mid \mathbf{y}, \mathbf{w}} \mathbf{f}(\mathbf{w}, z, \mathbf{y}) - E_{\mathbf{y}, z \mid \mathbf{w}} \mathbf{f}(\mathbf{w}, z, \mathbf{y}) \quad [11.28]$$

Both expectations can be computed via bottom-up algorithms like inside-outside. Alternatively, we can again maximize rather than marginalize over derivations for an approximate solution. In either case, the first term of the gradient requires us to identify derivations z that are compatible with the logical form \mathbf{y} . This can be done in a bottom-up dynamic programming algorithm, by having each cell in the table $t[i, j]$ include both the constituent types (e.g., NP, S) that can derive the span $w_{i:j-1}$, as well as the set of possible logical forms. The inclusion of logical forms means that the resulting table may be much larger than in syntactic parsing, where it is limited by the number of non-terminals in the grammar. This can be controlled by using pruning to eliminate intermediate analyses that are incompatible with the final logical form \mathbf{y} (Zettlemoyer and Collins, 2005), or by using beam search and restricting the size of each cell to some fixed constant (Liang et al., 2013).

If we replace each expectation in Equation 11.28 with argmax and then apply stochastic gradient descent to learn the weights, we obtain the **latent variable perceptron**, a simple and general algorithm for learning with missing data. The algorithm is shown in its most basic form in Algorithm 14, but the usual tricks such as averaging and margin loss can be applied (Yu and Joachims, 2009). Aside from semantic parsing, the latent variable perceptron has been used in tasks such as machine translation (Liang et al., 2006) and named entity recognition (Sun et al., 2009). In **latent conditional random fields**, we use the full expectations rather than maximizing over the hidden variable. This model has also been employed in a range of problems beyond semantic parsing, including parse reranking (Koo and Collins, 2005) and gesture recognition (Quattoni et al., 2007).

11.4.3 Learning from denotations

Logical forms are easier to obtain than complete derivations, but the annotation of logical forms still requires considerable expertise. However, it is relatively easy to obtain denotations for many natural language sentences. For example, in the geography domain, the

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 14 Latent variable perceptron

```

1: procedure LATENTVARIABLEPERCEPTRON( $\mathbf{w}^{(1:N)}, \mathbf{y}^{(1:N)}$ )
2:    $\theta \leftarrow \mathbf{0}$ 
3:   repeat
4:     Select an instance  $i$ 
5:      $\mathbf{z}^{(i)} \leftarrow \operatorname{argmax}_{\mathbf{z}} \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{z}, \mathbf{y}^{(i)})$ 
6:      $\hat{\mathbf{y}}, \hat{\mathbf{z}} \leftarrow \operatorname{argmax}_{\mathbf{y}', \mathbf{z}'} \theta \cdot \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{z}', \mathbf{y}')$ 
7:      $\theta \leftarrow \theta + \mathbf{f}(\mathbf{w}^{(i)}, \mathbf{z}^{(i)}, \mathbf{y}^{(i)}) - \mathbf{f}(\mathbf{w}^{(i)}, \hat{\mathbf{z}}, \hat{\mathbf{y}})$ 
8:   until tired
9:   return  $\theta$ 

```

denotation of a question would be its answer (Clarke et al., 2010; Liang et al., 2013):

Text : *What states border Georgia?*

Logical form : $\lambda x. \text{STATE}(x) \wedge \text{BORDER}(x, \text{GEORGIA})$

Denotation : {Alabama, Florida, North Carolina,
South Carolina, Tennessee}

Similarly, in a robotic control setting, the denotation of a command would be an action or sequence of actions (Artzi and Zettlemoyer, 2013). In both cases, the idea is to reward the semantic parser for choosing an analysis whose denotation is correct: the right answer to the question, or the right action.

Learning from logical forms was made possible by summing or maxing over derivations. This idea can be carried one step further, summing or maxing over all logical forms with the correct denotation. Let $v_i(\mathbf{y}) \in \{0, 1\}$ be a **validation function**, which assigns a binary score indicating whether the denotation $\llbracket \mathbf{y} \rrbracket$ for the text $\mathbf{w}^{(i)}$ is correct. We can then learn by maximizing a conditional-likelihood objective,

$$\ell^{(i)}(\theta) = \log \sum_{\mathbf{y}} v_i(\mathbf{y}) \times \mathbf{p}(\mathbf{y} \mid \mathbf{w}; \theta) \quad [11.29]$$

$$= \log \sum_{\mathbf{y}} v_i(\mathbf{y}) \times \sum_{\mathbf{z}} \mathbf{p}(\mathbf{y}, \mathbf{z} \mid \mathbf{w}; \theta), \quad [11.30]$$

which sums over all derivations \mathbf{z} of all valid logical forms, $\{\mathbf{y} : v_i(\mathbf{y}) = 1\}$. This corresponds to the log-probability that the semantic parser produces a logical form with a valid denotation.

Differentiating with respect to θ , we obtain,

$$\frac{\partial \ell^{(i)}}{\partial \theta} = \sum_{\mathbf{y}, \mathbf{z} : v_i(\mathbf{y})=1} \mathbf{p}(\mathbf{y}, \mathbf{z} \mid \mathbf{w}) \mathbf{f}(\mathbf{w}, \mathbf{z}, \mathbf{y}) - \sum_{\mathbf{y}', \mathbf{z}'} \mathbf{p}(\mathbf{y}', \mathbf{z}' \mid \mathbf{w}) \mathbf{f}(\mathbf{w}, \mathbf{z}', \mathbf{y}'), \quad [11.31]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

which is the usual difference in feature expectations. The positive term computes the expected feature expectations conditioned on the denotation being valid, while the second term computes the expected feature expectations according to the current model, without regard to the ground truth. Large-margin learning formulations are also possible for this problem. For example, Artzi and Zettlemoyer (2013) generate a set of valid and invalid derivations, and then impose a constraint that all valid derivations should score higher than all invalid derivations. This constraint drives a perceptron-like learning rule.

Additional resources

A key issue not considered here is how to handle **semantic underspecification**: cases in which there are multiple semantic interpretations for a single syntactic structure. Quantifier scope ambiguity is a classic example. Blackburn and Bos (2005) enumerate a number of approaches to this issue, and also provide links between natural language semantics and computational inference techniques. Much of the contemporary research on semantic parsing uses the framework of combinatory categorial grammar (CCG). Carpenter (1997) provides a comprehensive treatment of how CCG can support compositional semantic analysis. Another recent area of research is the semantics of multi-sentence texts. This can be handled with models of **dynamic semantics**, such as dynamic predicate logic (Groenendijk and Stokhof, 1991).

To learn more about ongoing research on data-driven semantic parsing, readers may consult the survey article by Liang and Potts (2015), tutorial slides and videos by Artzi and Zettlemoyer (2013),¹² and the source code by Yoav Artzi¹³ and Percy Liang.¹⁴

Exercises

1. Derive the **modus ponens** inference rule, which states that if we know $\phi \Rightarrow \psi$ and ϕ , then ψ must be true. The derivation can be performed using the definition of the \Rightarrow operator and some of the laws provided in § 11.2.1, plus one additional identity: $\perp \vee \phi = \phi$.
2. Convert the following examples into first-order logic, using the relations CAN-SLEEP, MAKES-NOISE, and BROTHER.
 - If Abigail makes noise, no one can sleep.
 - If Abigail makes noise, someone cannot sleep.

¹²Videos are currently available at <http://yoavartzi.com/tutorial/>

¹³<http://yoavartzi.com/spf>

¹⁴<https://github.com/percyliang/sempr>

- None of Abigail's brothers can sleep.
 - If one of Abigail's brothers makes noise, Abigail cannot sleep.
3. Extend the grammar fragment G_1 to include the ditransitive verb *teaches* and the proper noun *Swahili*. Show how to derive the interpretation for the sentence *Alex teaches Brit Swahili*, which should be $\text{TEACHES}(\text{ALEX}, \text{BRIT}, \text{SWAHILI})$. The grammar need not be in Chomsky Normal Form. For the ditransitive verb, use NP_1 and NP_2 to indicate the two direct objects.
 4. Derive the semantic interpretation for the sentence *Alex likes every dog*, using grammar fragment G_2 .
 5. Extend the grammar fragment G_2 to handle adjectives, so that the meaning of *an angry dog* is $\lambda P. \exists x \text{DOG}(x) \wedge \text{ANGRY}(x) \wedge P(x)$. Specifically, you should supply the lexical entry for the adjective *angry*, and you should specify the syntactic-semantic productions $\text{NP} \rightarrow \text{DET NOM}$, $\text{NOM} \rightarrow \text{JJ NOM}$, and $\text{NOM} \rightarrow \text{NN}$.
 6. Extend your answer to the previous question to cover copula constructions with predicative adjectives, such as *Alex is angry*. The interpretation should be $\text{ANGRY}(\text{ALEX})$. You should add a verb phrase production $\text{VP} \rightarrow \text{V}_{\text{cop}} \text{JJ}$, and a terminal production $\text{V}_{\text{cop}} \rightarrow \text{is}$. Show why your grammar extensions result in the correct interpretation.
 7. In Figure 11.6 and Figure 11.7, we treat the plurals *shoots* and *leaves* as entities. Revise G_2 so that the interpretation of *Alex eats leaves* is $\forall x. (\text{LEAF}(x) \Rightarrow \text{EATS}(\text{ALEX}, x))$, and show the resulting perceptron update.
 8. Statements like *every student eats a pizza* have two possible interpretations, depending on quantifier scope:

$$\forall x \exists y \text{PIZZA}(y) \wedge (\text{STUDENT}(x) \Rightarrow \text{EATS}(x, y)) \quad [11.32]$$

$$\exists y \forall x \text{PIZZA}(y) \wedge (\text{STUDENT}(x) \Rightarrow \text{EATS}(x, y)) \quad [11.33]$$

Explain why these interpretations really are different, and modify the grammar G_2 so that it can produce both interpretations.

9. Derive Equation 11.27.
10. **[todo: not sure this works]** Download the GeoQuery data, get some deterministic parser that overgenerates, and try to learn a reranker that selects the correct logical form.
11. In the GeoQuery domain, give a natural language query that has multiple plausible semantic interpretations with the same denotation. List both interpretations and the denotation.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Hint: There are many ways to do this, but one approach involves using toponyms (place names) that could plausibly map to several different entities in the model.

Chapter 12

Predicate-argument semantics

In this chapter, we consider more “lightweight” semantic representations. These semantic representations discard some aspects of first-order predicate calculus, but focus on predicate-argument structures. Let’s start with an example sentence:

(12.1) Asha gives Boyang a book.

The predicate calculus representation of this sentence would be written,

$$\exists x. \text{BOOK}(x) \wedge \text{GIVE}(\text{Asha}, \text{Boyang}, x) \quad [12.1]$$

In this representation, we define variable x for the book, and we link the strings *Asha* and *Boyang* to entities `Asha` and `Boyang`. Because the action of giving involves a giver, a recipient, and a gift, the predicate `GIVE` must take three arguments.

Now suppose we have additional information about the event, such as,

(12.2) Yesterday, Asha reluctantly gave Boyang a book.

One possible solution is to extend the predicate `GIVE` to take additional arguments,

$$\exists x. \text{BOOK}(x) \wedge \text{GIVE}(\text{Asha}, \text{Boyang}, x, \text{yesterday}, \text{reluctantly}) \quad [12.2]$$

But this is clearly unsatisfactory: *yesterday* and *reluctantly* are optional arguments, and we would need a different version of the `GIVE` predicate for every possible combination of arguments. **Event semantics** solves this problem by **reifying** the event as an existentially quantified variable e ,

$$\begin{aligned} \exists e, x. & \text{GIVE-EVENT}(e) \wedge \text{GIVER}(e, \text{Asha}) \wedge \text{GIFT}(e, x) \wedge \text{BOOK}(e, x) \wedge \text{RECIPIENT}(e, \text{Boyang}) \\ & \wedge \text{TIME}(e, \text{Yesterday}) \wedge \text{MANNER}(e, \text{reluctantly}) \end{aligned}$$

In this way, each argument of the event — the giver, the recipient, the gift — can be represented with a relation of its own, linking the argument to the event e . The expression $\text{GIVER}(e, \text{Asha})$ says that Asha plays the **role** of GIVER in the event. This reformulation nicely handles the problem of optional information such as the time or manner of the event, which are called **adjuncts**. Unlike arguments, adjuncts are not a mandatory part of the relation, but under this representation, they can be expressed with additional logical relations that are conjoined to the semantic interpretation of the sentence.¹

The event semantic representation can be applied to nested clauses, e.g.,

(12.3) Chris sees Asha pay Boyang.

This is done by using the event variable as an argument:

$$\begin{aligned} \exists e_1, e_2. & \text{SEE-EVENT}(e_1) \wedge \text{SEER}(e_1, \text{Chris}) \wedge \text{SIGHT}(e_1, e_2) \\ & \wedge \text{PAY-EVENT}(e_2) \wedge \text{PAYER}(e_2, \text{Asha}) \wedge \text{PAYEE}(e_2, \text{Boyang}) \end{aligned} \quad [12.3]$$

As with first-order predicate calculus, the goal of event semantics is to provide a representation that generalizes over many surface forms. Consider the following paraphrases of (12.1):

(12.4) Asha gives a book to Boyang.

(12.5) A book is given to Boyang by Asha.

(12.6) A book is given by Asha to Boyang.

(12.7) The gift of a book from Asha to Boyang ...

All have the same event semantic meaning, given in Equation 12.1. Note that the final example does not include a verb! Events are often introduced by verbs, but not always: in this final example, the noun *gift* introduces the same predicate, with the same accompanying arguments.

Semantic role labeling (SRL) is a relaxed form of semantic parsing, in which each semantic role is filled by a set of tokens from the text itself. This is sometimes called “shallow semantics” because, unlike model-theoretic semantic parsing, role fillers need not be symbolic expressions with denotations in some world model. A semantic role labeling system is required to identify all predicates, and then specify the spans of text that fill each role, when possible. To get a sense of the task, here is a more complicated example:

¹This representation is often called **Neo-Davidsonian event semantics**. The use of existentially-quantified event variables was proposed by Davidson (1967) to handle the issue of optional adjuncts. In Neo-Davidsonian semantics, this treatment of adjuncts is extended to mandatory arguments as well (e.g., Parsons, 1990).

(12.8) Boyang wants Asha to give him a linguistics book.

In this example, there are two predicates, expressed by the verbs *want* and *give*. Thus, a semantic role labeler should return the following output:

- (PREDICATE : *wants*, WANTER : *Boyang*, DESIRE : *Asha to give him a linguistics book*)
- (PREDICATE : *give*, GIVER : *Asha*, RECIPIENT : *him*, GIFT : *a linguistics book*)

In the example, *Boyang* and *him* may refer to the same person, but this would be ignored in semantic role labeling. However, in other predicate-argument representations, such as **Abstract Meaning Representation (AMR)**, such references must be resolved. We will return to AMR in § 12.3, but first, let us further consider the notion of semantic roles.

12.1 Semantic roles

As discussed so far, event semantics requires specifying a number of additional logical relations to link arguments to events: GIVER, RECIPIENT, SEER, SIGHT, etc. Indeed, every predicate requires a set of logical relations to express its own arguments. In contrast, adjuncts such as TIME and MANNER are shared across many types of events. A natural question is whether it is possible to treat mandatory arguments more like adjuncts, by identifying a set of generic argument types that are shared across many event predicates. This can be further motivated by examples involving semantically related verbs:

(12.9) Asha gave Boyang a book.

(12.10) Asha loaned Boyang a book.

(12.11) Asha taught Boyang a lesson.

(12.12) Asha gave Boyang a lesson.

In the first two examples, the roles of Asha, Boyang, and the book are nearly identical. The third example is slightly different, but the fourth example shows that the roles of GIVER and TEACHER can be viewed as related.

One way to think about the relationship between roles such as GIVER and TEACHER is by enumerating the set of properties that an entity typically possesses when it fulfills these roles: givers and teachers are usually animate and “volitional” (meaning that they choose to enter into the action).² In contrast, the thing that gets loaned or taught is usually not animate or volitional; furthermore, it is unchanged by the event.

²There are always exceptions. For example, in the sentence *The C programming language has taught me a lot about perseverance*, the “teacher” is the *The C programming language*, which is presumably not animate or volitional.

	<i>Asha</i>	<i>gave</i>	<i>Boyang</i>	<i>a book</i>
VerbNet	AGENT		RECIPIENT	THEME
PropBank	ARG0: giver		ARG2: entity given to	ARG1: thing given
FrameNet	DONOR		RECIPIENT	THEME
	<i>Asha</i>	<i>taught</i>	<i>Boyang</i>	<i>algebra</i>
VerbNet	AGENT		RECIPIENT	TOPIC
PropBank	ARG0: teacher		ARG2: student	ARG1: subject
FrameNet	TEACHER		STUDENT	SUBJECT

Figure 12.1: Example semantic annotations according to VerbNet, PropBank, and FrameNet

Building on these ideas, **thematic roles** generalize across predicates by leveraging the shared semantic properties of typical role fillers (Fillmore, 1968). For example, in examples (12.9-12.12), Asha plays a similar role in all four sentences, which we will call the **agent**. This reflects a number of shared semantic properties: she is the one who is actively and intentionally performing the action, while Boyang is a more passive participant; the book and the lesson would play a different role, as non-animate participants in the event.

Let us now consider a few well-known approaches to semantic roles. Example annotations from each of these systems are shown in Figure 12.1.

12.1.1 VerbNet

VerbNet (Kipper-Schuler, 2005) is a lexicon of verbs, and it includes thirty “core” thematic roles played by arguments to these verbs. Here are some example roles, accompanied by their definitions from the VerbNet Guidelines.³

- **AGENT**: “ACTOR in an event who initiates and carries out the event intentionally or consciously, and who exists independently of the event.”
- **PATIENT**: “UNDERGOER in an event that experiences a change of state, location or condition, that is causally involved or directly affected by other participants, and exists independently of the event.”
- **RECIPIENT**: “DESTINATION that is animate”
- **THEME**: “UNDERGOER that is central to an event or state that does not have control over the way the event occurs, is not structurally changed by the event, and/or is characterized as being in a certain position or condition throughout the state.”

³http://verbs.colorado.edu/verb-index/VerbNet_Guidelines.pdf

- TOPIC: “THEME characterized by information content transferred to another participant.”

VerbNet roles are organized in a hierarchy, so that a TOPIC is a type of THEME, which in turn is a type of UNDERGOER, which is a type of PARTICIPANT, the top-level category.

In addition, VerbNet organizes verb senses into a class hierarchy, in which verb senses that have similar meanings are grouped together. Recall from § 3.2 that multiple meanings of the same word are called **senses**, and that WordNet identifies senses for many English words. VerbNet builds on WordNet, so that verb classes are identified by the WordNet senses of the verbs that they contain. For example, the verb class `give-13.1` includes the first WordNet sense of *loan* and the second WordNet sense of *lend*.

Each VerbNet class or subclass takes a set of thematic roles. For example, `give-13.1` takes arguments with the thematic roles of AGENT, THEME, and RECIPIENT;⁴ the predicate TEACH takes arguments with the thematic roles AGENT, TOPIC, RECIPIENT, and SOURCE.⁵ So according to VerbNet, *Asha* and *Boyang* play the roles of AGENT and RECIPIENT in the sentences,

(12.13) *Asha* gave *Boyang* a book.

(12.14) *Asha* taught *Boyang* algebra.

The *book* and *algebra* are both THEMES, but *algebra* is a subcategory of THEME — a TOPIC — because it consists of information content that is given to the receiver.

12.1.2 Proto-roles and PropBank

Detailed thematic role inventories of the sort used in VerbNet are not universally accepted. For example, (Dowty, 1991, pp. 547) notes that “Linguists have often found it hard to agree on, and to motivate, the location of the boundary between role types.” He argues that a solid distinction can be identified between just two **proto-roles**, which have a number of distinguishing characteristics:

- PROTO-AGENT: volitional involvement in the event or state; sentience and/or perception; causing an event or change of state in another participant; movement; exists independently of the event.

⁴<https://verbs.colorado.edu/verb-index/vn/give-13.1.php>

⁵https://verbs.colorado.edu/verb-index/vn/transfer_mesg-37.1.1.1.php

- **PROTO-PATIENT:** undergoes change of state; causally affected by another participant; stationary relative to the movement of another participant; does not exist independently of the event.⁶

In the examples in Figure 12.1, Asha has most of the proto-agent properties: in giving the book to Boyang, she is acting volitionally (as opposed to *Boyang got a book from Asha*, in which it is not clear whether Asha gave up the book willingly); she is sentient; she causes a change of state in Boyang; she exists independently of the event. Boyang has some proto-agent properties (for example, he is sentient and exists independently of the event), but he also some proto-patient properties: he is the one who is causally affected and who undergoes change of state in both cases. The book that Asha gives Boyang has fewer still of the proto-agent properties — it is not volitional or sentient, and it has no causal role — but it also has few proto-patient properties, as it does not undergo change of state and is not stationary.

The **Proposition Bank**, or PropBank (Palmer et al., 2005), builds on this basic agent-patient distinction, as a middle ground between generic thematic roles and predicate-specific “deep roles.” Each verb is linked to a list of numbered arguments, with ARG0 as the proto-agent and ARG1 as the proto-patient. Additional numbered arguments are verb-specific. For example, for the predicate TEACH,⁷ the arguments are:

- ARG0: the teacher
- ARG1: the subject
- ARG2: the student(s)

Verbs may have any number of arguments: for example, WANT and GET have five, while EAT has only ARG0 and ARG1. In addition to the semantic arguments found in the frame files, roughly a dozen general-purpose **adjuncts** may be used in combination with any verb. These are shown in Table 12.1.

PropBank-style semantic role labeling is annotated over the entire Penn Treebank. This annotation includes the sense of each verbal predicate, as well as the argument spans.

12.1.3 FrameNet

Semantic **frames** are descriptions of situations or events. Frames may be **evoked** by one of their **lexical units** (often a verb, but not always), and they include some number of

⁶Reisinger et al. (2015) ask crowd workers to annotate these properties directly, finding that annotators tend to agree on the properties of each argument. They also find that in English, arguments having more proto-agent properties tend to appear in subject position, while arguments with more proto-patient properties appear in object position.

⁷<http://verbs.colorado.edu/propbank/framesets-english-aliases/teach.html>

TMP	time	<i>Boyang ate a bagel</i> [AM-TMP <i>yesterday</i>].
LOC	location	<i>Asha studies in</i> [AM-LOC <i>Stuttgart</i>]
MOD	modal verb	<i>Asha</i> [AM-MOD <i>will</i>] <i>study in Stuttgart</i>
ADV	general purpose	[AM-ADV <i>Luckily</i>], <i>Asha knew algebra</i> .
MNR	manner	<i>Asha ate</i> [AM-MNR <i>aggressively</i>].
DIS	discourse connective	[AM-DIS <i>However</i>], <i>Asha prefers algebra</i> .
PRP	purpose	<i>Barry studied</i> [AM-PRP <i>to pass the bar</i>].
DIR	direction	<i>Workers dumped burlap sacks</i> [AM-DIR <i>into a bin</i>].
NEG	negation	<i>Asha does</i> [AM-NEG <i>not</i>] <i>know algebra</i> .
EXT	extent	<i>Prices increased</i> [AM-EXT <i>4%</i>].
CAU	cause	<i>Asha returned the book</i> [AM-CAU <i>because it was overdue</i>].

Table 12.1: PropBank adjuncts (Palmer et al., 2005), sorted by frequency in the corpus

frame elements, which are like roles (Fillmore, 1976). For example, the act of teaching is a frame, and can be evoked by the verb *taught*; the associated frame elements include the teacher, the student(s), and the subject being taught. Frame semantics has played a significant role in the history of artificial intelligence, in the work of Minsky (1974) and Schank and Abelson (1977). In natural language processing, the theory of frame semantics has been implemented in **FrameNet** (Fillmore and Baker, 2009), which consists of a lexicon of roughly 1000 frames, and a corpus of more than 200,000 “exemplar sentences,” in which the frames and their elements are annotated.⁸

Rather than seeking to link semantic roles such as TEACHER and GIVER into thematic roles such as AGENT, FrameNet aggressively groups verbs into frames, and links semantically-related roles across frames. For example, the following two sentences would be annotated identically in FrameNet:

(12.15) Asha taught Boyang algebra.

(12.16) Boyang learned algebra from Asha.

This is because *teach* and *learn* are both lexical units in the EDUCATION.TEACHING frame. Furthermore, roles can be shared even when the frames are distinct, as in the following two examples:

(12.17) Asha gave Boyang a book.

(12.18) Boyang got a book from Asha.

⁸These statistics are accurate at the time of this writing, in 2017. Current details can be found at the website, <https://framenet.icsi.berkeley.edu/>

The GIVING and GETTING frames both have RECIPIENT and THEME elements, so Boyang and the book would play the same role. Asha’s role is different: she is the DONOR in the GIVING frame, and the SOURCE in the GETTING frame. FrameNet makes extensive use of multiple inheritance to share information across frames and frame elements: for example, the COMMERCE_SELL and LENDING frames inherit from GIVING frame.

12.2 Semantic role labeling

The task of semantic role labeling is to identify the parts of the sentence comprising the semantic roles. In English, this task is typically performed on the PropBank corpus, with the goal of producing outputs in the following form:

(12.19) [ARG0 Asha] [GIVE.01 gave] [ARG2 Boyang’s mom] [ARG1 a book] [AM-TMP yesterday].

Note that a single sentence may have multiple verbs, and therefore a given word may be part of multiple role-fillers:

(12.20) [ARG0 Asha] [WANT.01 wanted] [ARG1 Boyang to give her the book].
 Asha wanted [ARG0 Boyang] [GIVE.01 to give] [ARG2 her] [ARG1 the book].

12.2.1 Semantic role labeling as classification

PropBank is annotated on the Penn Treebank, and annotators used phrasal constituents (§ 8.2.2) to fill the roles. Therefore SRL can be viewed as the task of assigning to each phrase a label from the set $\mathcal{R} = \{\emptyset, \text{PRED}, \text{ARG0}, \text{ARG1}, \text{ARG2}, \dots, \text{AM-LOC}, \text{AM-TMP}, \dots\}$, where \emptyset indicates that the phrase plays no role, and PRED indicates that it is the verbal predicate. If we treat semantic role labeling as a classification problem, we obtain the following functional form:

$$\hat{y}_{(i,j)} = \underset{y}{\operatorname{argmax}} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y, i, j, \rho, \tau), \quad [12.4]$$

where,

- (i, j) indicates the span of a phrasal constituent $(w_i, w_{i+1}, \dots, w_{j-1})$;⁹
- \mathbf{w} represents the sentence as a sequence of tokens;

⁹PropBank roles can also be filled by **split constituents**, which are discontinuous spans of text. This situation most frequently in reported speech, e.g. [ARG1 *By addressing these problems*], *Mr. Maxwell said*, [ARG1 *the new funds have become extremely attractive.*] (example adapted from Palmer et al., 2005). This issue is typically addressed by defining “continuation arguments”, e.g. C-ARG1, which refers to the continuation of ARG1 after the split.

Predicate lemma and POS tag	The lemma of the predicate verb and its part-of-speech tag
Voice	Whether the predicate is in active or passive voice, as determined by a set of syntactic patterns for identifying passive voice constructions
Phrase type	The constituent phrase type for the proposed argument in the parse tree, e.g. NP, PP
Headword and POS tag	The head word of the proposed argument and its POS tag, identified using the Collins (1997) rules
Position	Whether the proposed argument comes before or after the predicate in the sentence
Syntactic path	The set of steps on the parse tree from the proposed argument to the predicate (described in detail in the text)
Subcategorization	The syntactic production from the first branching node above the predicate. For example, in Figure 12.2, the subcategorization feature around <i>taught</i> would be VP \rightarrow VBD NP PP.

Table 12.2: Features used in semantic role labeling by Gildea and Jurafsky (2002).

- ρ is the index of the predicate verb in w ;
- τ is the structure of the phrasal constituent parse of w .

Table 12.2 shows the features used in the seminal paper on FrameNet semantic role labeling by Gildea and Jurafsky (2002). By 2005 there were several systems for PropBank semantic role labeling, and their approaches and feature sets are summarized by Carreras and Màrquez (2005). Typical features include: the phrase type, head word, part-of-speech, boundaries, and neighbors of the proposed argument $w_{i:j}$; the word, lemma, part-of-speech, and voice of the verb w_ρ (active or passive), as well as features relating to its frameset; the distance and path between the verb and the proposed argument. In this way, semantic role labeling systems are high-level “consumers” in the NLP stack, using features produced from lower-level components such as part-of-speech taggers and parsers. More comprehensive and contemporary feature sets are enumerated by Das et al. (2014) and Täckström et al. (2015).

A particularly powerful class of features relate to the **syntactic path** between the argument and the predicate. These features capture the sequence of moves required to get from the argument to the verb by traversing the phrasal constituent parse of the sentence. The idea of these features is to capture syntactic regularities in how various arguments are realized. Syntactic path features are best illustrated by example, using the parse tree

(c) Jacob Eisenstein 2014-2017. Work in progress.

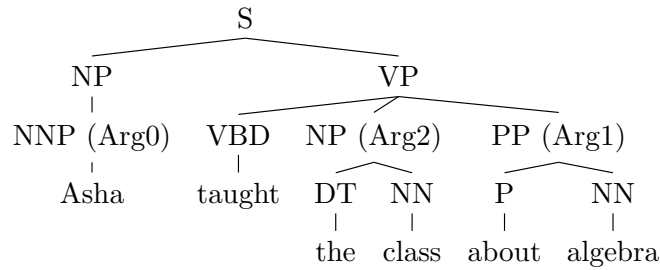


Figure 12.2: Semantic role labeling is often performed on the parse tree for a sentence, labeling individual constituents. [todo: check arg1; show arrows for path features]

in Figure 12.2:

- The path from *Asha* to the verb *taught* is $\text{NNP}\uparrow\text{NP}\uparrow\text{S}\downarrow\text{VP}\downarrow\text{VBD}$. The first part of the path, $\text{NNP}\uparrow\text{NP}\uparrow\text{S}$, means that we must travel up the parse tree from the NNP tag (proper noun) to the S (sentence) constituent. The second part of the path, $\text{S}\downarrow\text{VP}\downarrow\text{VBD}$, means that we reach the verb by producing a VP (verb phrase) from the S constituent, and then by producing a VBD (past tense verb). This feature is consistent with *Asha* being in subject position, since the path includes the sentence root S.
- The path from *the class* to the verb is $\text{NP}\uparrow\text{VP}\downarrow\text{VBD}$. This is consistent with *the class* being in object position, since the path passes through the VP node that dominates the verb *taught*.

Because there are many possible path features, it can also be helpful to look at smaller parts: for example, the upward and downward parts can be treated as separate features; another feature might consider whether S appears anywhere in the path.

Rather than using the constituent parse, it is also possible to build features from the **dependency path** between the head word of each argument and the verb (Pradhan et al., 2005). Using the Universal Dependency part-of-speech tagset and dependency relations (Nivre et al., 2016), the dependency path from *Asha* to *taught* is $\text{PROPN} \xleftarrow{\text{NSUBJ}} \text{VERB}$, because *taught* is the head of a relation of type $\xleftarrow{\text{NSUBJ}}$ with *Asha*. Similarly, the dependency path from *class* to *taught* is $\text{NOUN} \xleftarrow{\text{DOBJ}} \text{VERB}$, because *class* heads the noun phrase that is a direct object of *taught*. A more interesting example is *Asha tried to teach the class*, where the path from *Asha* to *tried* is $\text{PROPN} \xleftarrow{\text{NSUBJ}} \text{VERB} \rightarrow \text{VERB}$. The right-facing arrow in second relation indicates that *tried* is the head of its XCOMP relation with *teach*.

(c) Jacob Eisenstein 2014-2017. Work in progress.

12.2.2 Semantic role labeling as constrained optimization

A potential problem with treating SRL as a classification problem is that there are a number of sentence-level **constraints**, which a classifier might violate.

- For a given verb, there can be only one argument of each type (ARG0, ARG1, etc.)
- Arguments cannot overlap. This problem arises when we are labeling the phrases in a constituent parse tree, as shown in Figure 12.2: if we label the PP *about algebra* as an argument or adjunct, then its children *about* and *algebra* must be labeled as \emptyset . The same constraint also applies to the syntactic ancestors of this phrase.

These constraints can be viewed as introducing dependencies across labeling decisions. In structure prediction problems such as sequence labeling and parsing, such dependencies are usually handled by defining additional features over the entire structure, y . Efficient inference requires that the global features have a local decomposition that enables dynamic programming: for example, in sequence labeling, the features over y were decomposed into features over pairs of adjacent tags, permitting the application of the Viterbi algorithm for inference. Unfortunately, the constraints that arise in semantic role labeling are less amenable to local decomposition — particularly the constraint that each argument is used only once in the sentence.¹⁰ We therefore consider **constrained optimization** as an alternative solution.

Let the **feasible set** $\mathcal{C}(\tau)$ refer to all labelings that obey the constraints introduced by the parse τ . We can reformulate the semantic role labeling problem as a constrained optimization,

$$\begin{aligned} \max_{\mathbf{y}} \quad & \sum_{(i,j) \in \tau} \theta \cdot \mathbf{f}(\mathbf{w}, y_{(i,j)}, i, j, \rho, \tau) \\ \text{s.t.} \quad & \mathbf{y} \in \mathcal{C}(\tau). \end{aligned} \tag{12.5}$$

In this formulation, the objective (shown on the first line) is a separable function of each individual labeling decision, but the constraints (shown on the second line) apply to the overall labeling. The sum $\sum_{(i,j) \in \tau}$ indicates that we are summing over all constituent spans in the parse τ . The expression *s.t.* in the second line means that we maximize the objective *subject to* the constraints that appear to the right.

Integer linear programming

A number of practical algorithms exist for restricted forms of constrained optimization. One such restricted form is **integer linear programming**, in which we optimize a linear

¹⁰Dynamic programming solutions have been proposed by Tromble and Eisner (2006) and Täckström et al. (2015), but they involve creating a trellis structure whose size is exponential in the number of labels.

objective function over integer variables, with linear constraints. To formulate SRL as an integer linear program, we begin by rewriting the labels as a set of binary variables $\mathbf{z} = \{z_{i,j,r}\}$, where,

$$z_{i,j,r} = \begin{cases} 1, & y_{(i,j)} = r \\ 0, & \text{otherwise.} \end{cases} \quad [12.6]$$

Thus, the variables \mathbf{z} are a binarized version of the semantic role labeling \mathbf{y} .

Objective Next, we restrict the objective to be a linear function of \mathbf{z} . We begin with the feature function, $\mathbf{f}(\mathbf{w}, y_{(i,j)}, i, j, \rho, \tau)$. Such features are typically logical conjunctions involving the label $y_{(i,j)}$. For example:

$$f_j(\mathbf{w}, y_{(i,j)}, i, j, \rho, \tau) = \begin{cases} 1, & y_{(i,j)} = \text{ARG1} \wedge w_i = \textit{the} \\ 0, & \text{otherwise.} \end{cases} \quad [12.7]$$

This feature is an indicator that takes the value 1 for constituents that are labeled ARG1 and begin with the word *the*. If all features are conjunctions with the label, then the feature function can be rewritten,

$$\mathbf{f}(\mathbf{w}, y_{(i,j)}, i, j, \rho, \tau) = \sum_{r \in \mathcal{R}} z_{i,j,r} \times \mathbf{g}(\mathbf{w}, i, j, r, \rho, \tau), \quad [12.8]$$

where \mathcal{R} is the set of all possible labels, $\{\text{ARG0}, \text{ARG1}, \dots, \text{AM-LOC}, \dots, \emptyset\}$. With this change in notation, we can now rewrite the objective,

$$\sum_{(i,j) \in \tau} \boldsymbol{\theta} \cdot \mathbf{f}(\mathbf{w}, y_{(i,j)}, i, j, \rho, \tau) = \sum_{(i,j) \in \tau} \boldsymbol{\theta} \cdot \left(\sum_{r \in \mathcal{R}} z_{i,j,r} \times \mathbf{g}(\mathbf{w}, i, j, r, \rho, \tau) \right) \quad [12.9]$$

$$= \sum_{(i,j) \in \tau} \sum_{r \in \mathcal{R}} z_{i,j,r} (\boldsymbol{\theta} \cdot \mathbf{g}(\mathbf{w}, i, j, r, \rho, \tau)) \quad [12.10]$$

$$= \sum_{(i,j) \in \tau} \sum_{r \in \mathcal{R}} z_{i,j,r} \psi_{i,j,r}, \quad [12.11]$$

where $\psi_{i,j,r} \triangleq \boldsymbol{\theta} \cdot \mathbf{g}(\mathbf{w}, i, j, r, \rho, \tau)$. This objective is clearly a linear function of the variables $\mathbf{z} = \{z_{i,j,r}\}$.

Constraints Integer linear programming permits linear inequality constraints, of the general form $\mathbf{A}\mathbf{z} \leq \mathbf{b}$, where the parameters \mathbf{A} and \mathbf{b} define the constraints. To make this more concrete, let's start with the constraint that each non-null role type can occur only once in a sentence. This constraint can be written,

$$\forall r \neq \emptyset, \quad \sum_{(i,j) \in \tau} z_{i,j,r} \leq 1. \quad [12.12]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

Recall that $z_{i,j,r} = 1$ if and only if the span (i, j) has label r ; this constraint says that for each possible label $r \neq \emptyset$, there can be at most one (i, j) such that $z_{i,j,r} = 1$. This constraint can be written in the form $\mathbf{Az} \leq \mathbf{b}$, as you will find if you complete the exercises at the end of the chapter.

Now consider the constraint that labels cannot overlap. Let the function $\pi_\tau(i, j) = \{(i', j')\}$ indicate the set of constituents that are ancestors or descendants of (i, j) in the parse τ . For any (i, j) such that $y_{i,j} \neq \emptyset$, the non-overlapping constraint means that all of its ancestors and descendants (i', j') must be labeled as a non-argument, $y_{i',j'} = \emptyset$. We can write this as a set of linear constraints,

$$\forall (i, j) \in \tau, \quad \sum_{r \neq \emptyset} \left(z_{i,j,r} + \sum_{(i',j') \in \pi_\tau(i,j)} z_{i',j',r} \right) \leq 1. \quad [12.13]$$

We can therefore rewrite the semantic role labeling problem as the following integer linear program,

$$\max_{\mathbf{z} \in \{0,1\}^{|\tau|}} \sum_{(i,j) \in \tau} \sum_{r \in \mathcal{R}} z_{i,j,r} \psi_{i,j,r} \quad [12.14]$$

$$\text{s.t.} \quad \forall r \neq \emptyset, \quad \sum_{(i,j) \in \tau} z_{i,j,r} \leq 1. \quad [12.15]$$

$$\forall (i, j) \in \tau, \quad \sum_{r \neq \emptyset} \left(z_{i,j,r} + \sum_{(i',j') \in \pi_\tau(i,j)} z_{i',j',r} \right) \leq 1. \quad [12.16]$$

The effectiveness of integer linear programming for semantic role labeling was first demonstrated by Punyakanok et al. (2008).

Learning with constraints Learning can be performed in the context of constrained optimization using the usual perceptron or large-margin classification updates. Because constrained inference is generally more time-consuming, a key question is whether it is necessary to apply the constraints during learning. Chang et al. (2008) find that better performance can be obtained by learning *without* constraints, and then applying constraints only when using the trained model to predict semantic roles for unseen data.

How important are the constraints? Das et al. (2014) find that an unconstrained, classification-based method performs nearly as well as constrained optimization for FrameNet parsing: while it commits many violations of the “no-overlap” constraint, the overall F_1 score is less than one point worse than the score at the constrained optimum. Similar results are obtained for PropBank semantic role labeling by Punyakanok et al. (2008). He et al.

(c) Jacob Eisenstein 2014-2017. Work in progress.

(2017) find that constrained inference makes a bigger impact if the constraints are based on manually-labeled “gold” syntactic parses. This implies that errors from the syntactic parser may limit the effectiveness of the constraints. Punyakanok et al. (2008) hedge against parser error by including constituents from several different parsers; any constituent can be selected from any parse, and additional constraints ensure that overlapping constituents are not selected.

Implementation Integer linear programming solvers such as `glpk`,¹¹ `cplex`,¹² and `Gurobi`¹³ allow inequality constraints to be expressed directly in the problem definition, rather than in the matrix form $\mathbf{Az} \leq \mathbf{b}$. The time complexity of integer linear programming is theoretically exponential in the number of variables $|z|$, but in practice these off-the-shelf solvers obtain good solutions efficiently. Das et al. (2014) report that the `cplex` solver requires 43 seconds to perform inference on the FrameNet test set, which contains 4,458 predicates.

Recent work has shown that many constrained optimization problems in natural language processing can be solved in a highly parallelized fashion, using optimization techniques such as **dual decomposition**, which are capable of exploiting the underlying problem structure (Rush et al., 2010). Das et al. (2014) apply this technique to FrameNet semantic role labeling, obtaining an order-of-magnitude speedup over `cplex`.

12.2.3 Neural semantic role labeling

Neural network approaches to SRL have tended to treat it as a sequence labeling task, using a labeling scheme such as the **BIO notation**, which we previously saw in named entity recognition (§ 7.3). In this notation, the first token in a span of type ARG1 is labeled B-ARG1; all remaining tokens in the span are **inside**, and are therefore labeled I-ARG1. Tokens outside any argument are labeled O. For example:

(12.21) *Asha taught Boyang 's mom about algebra*
 B-ARG0 PRED B-ARG2 I-ARG2 I-ARG2 B-ARG1 I-ARG1

We now consider two classes of neural networks that can learn to produce such labelings.

Convolutional neural networks One of the first applications of **convolutional neural networks** (§ 2.4) to natural language processing was the task of classifying semantic roles. Collobert et al. (2011b) treat the task as a classification problem, using information gathered from across the sentence to compute the label for each token. For example, suppose our goal is to tag the role of word m with respect to verb v ; then for word n , we compute the discrete feature vector, $\mathbf{f}(\mathbf{w}, n, v, m)$, which would include the lower-case word

¹¹<https://www.gnu.org/software/glpk/>

¹²<https://www-01.ibm.com/software/commerce/optimization/cplex-optimizer/>

¹³<http://www.gurobi.com/>

w_m , and the distances $m - n$ and $m - v$. These features are then used as the inputs to a nonlinear prediction model, as follows:

- Each discrete feature is associated with a dense vector embedding, and these embeddings are concatenated, resulting in a dense vector $\mathbf{x}_m^{(0)}$. The horizontal concatenation of the dense embeddings for all words in the sentence is $\mathbf{X}^{(0)} = [\mathbf{x}_0^{(0)}, \mathbf{x}_1^{(0)}, \dots, \mathbf{x}_M^{(0)}]$.
- Next, a convolutional operation is applied to merge information across words, $\mathbf{X}^{(1)} = \mathbf{C}\mathbf{X}^{(0)}$. Thus, $\mathbf{x}_m^{(1)}$ contains information about the word w_m , but also about its near neighbors. (Special padding vectors are included on the left and right ends of the matrix $\mathbf{X}^{(0)}$ before convolution.)
- To convert the matrix $\mathbf{X}^{(1)}$ back to a vector $\mathbf{z}^{(1)}$, Collobert *et al.* apply **max pooling**. We will write $\mathbf{z} = \text{MaxPool}(\mathbf{X})$ to indicate that each $z_j = \max_m(x_{0,j}^{(1)}, x_{1,j}^{(1)}, \dots, x_{M,j}^{(1)})$.
- The vector $\mathbf{z}^{(1)}$ is then passed through several feedforward layers, $\mathbf{z}^{(i)} = \mathbf{g}(\mathbf{\Theta}^{(i)} \mathbf{z}^{(i-1)})$, where $\mathbf{\Theta}^{(i)}$ is a matrix of weights and \mathbf{g} is an elementwise nonlinear transformation.¹⁴
- At the output layer $\mathbf{z}^{(K)}$ is used to make a prediction, $\hat{y} = \arg\max_y \mathbf{\Theta}^{(y)} \mathbf{z}^{(K)}$. [todo: consider making this a figure/algorithm]

Collobert et al. (2011b) apply this tagging model without regard to constraints, so in principle it could produce labelings that include each argument multiple times. The parameters of the model include the word and feature embeddings that constitute $\mathbf{X}^{(0)}$, the convolution matrix \mathbf{C} , the feedforward weight matrices $\mathbf{\Theta}^{(i)}$, and the prediction weights $\mathbf{\Theta}^{(y)}$. Each of these parameters is estimated by backpropagated stochastic gradient descent (see § 5.3.1). Collobert et al. (2011b) emphasize that **multi-task learning** was essential to get good performance: they train the word embeddings not only to accurately predict PropBank labels, but also to assign a high likelihood to a large corpus of unlabeled data. A more contemporary approach would be to use **pre-trained word embeddings**, which have already been trained to predict words in context, thereby avoiding the cost of jointly training across a large unlabeled dataset.

¹⁴Collobert *et al.* use a piecewise linear **hard tanh** function for nonlinear transformations,

$$g(x) = \begin{cases} -1, & x < -1 \\ x, & -1 \leq x \leq 1 \\ 1, & x > 1. \end{cases} \quad [12.17]$$

An advantage of this function is that the gradient is easy and fast to compute, making training faster. More recent work has emphasized the **rectified linear unit (ReLU)**, $g(x) = \max(x, 0)$. This function, which is described in chapter 5, also has similar advantages to hard tanh, but avoids saturation because it has a non-zero gradient for large values of x (Salinas and Abbott, 1996; Glorot et al., 2011).

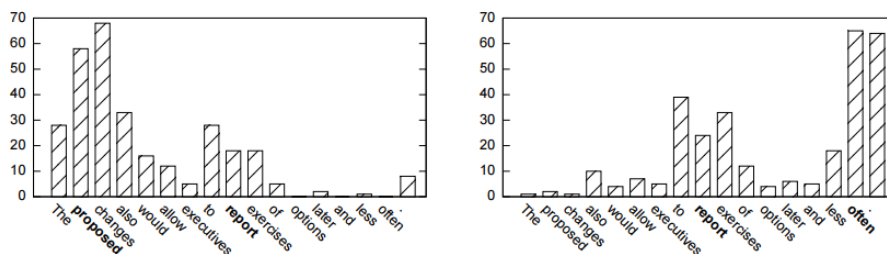


Figure 12.3: Number of features chosen at each word position by the max pooling operation, for tagging the words *proposed* (left) and *often* (right). Figure reprinted from Collobert et al. (2011b) [todo: ask for permission]

A key aspect of convolutional neural networks is the use of **pooling** operations, which combine information across a variable-length sequence of vectors into a single vector or matrix. Max pooling is widely used in natural language processing applications, because it enables each element in the vector \mathbf{z} to take information from across a sentence or other sequence of text. Figure 12.3 shows the number of “features” taken from each word in a sentence — that is, how often the max operation chooses an element from each word. In each case, the pooling operation emphasizes the word to be tagged, its neighbors, and also the main verb *report*.

Recurrent neural networks An alternative neural approach to semantic role labeling is to use **recurrent neural network** models, such as **long short-term memories** (LSTMs; see § 6.5.4 to review how these models are applied to tagging tasks). Zhou and Xu (2015) apply a bidirectional multilayer LSTM to PropBank semantic role labeling. In this model, each bidirectional LSTM serves as input for another, higher-level bidirectional LSTM, allowing complex non-linear transformations of the original input embeddings, $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_M]$. The hidden state of the final LSTM is $\mathbf{Z}^{(K)} = [z_1^{(K)}, z_2^{(K)}, \dots, z_M^{(K)}]$. The “emission” score for each tag $Y_m = y$ is equal to the inner product $\theta_y \cdot z_m^{(K)}$, and there is also a transition score for each pair of adjacent tags. The complete model can be written,

$$\mathbf{Z}^{(1)} = \text{BiLSTM}(\mathbf{X}) \quad [12.18]$$

$$\mathbf{Z}^{(i)} = \text{BiLSTM}(\mathbf{Z}^{(i-1)}) \quad [12.19]$$

$$\hat{\mathbf{y}} = \underset{\mathbf{y}}{\operatorname{argmax}} \sum_{m=1}^M \Theta^{(y)} z_m^{(K)} + \psi_{y_{m-1}, y_m}. \quad [12.20]$$

Note that the final step maximizes over the entire labeling \mathbf{y} , and includes a score for each tag transition ψ_{y_{m-1}, y_m} . This combination of LSTM and pairwise potentials on tags is an example of an **LSTM-CRF**. The maximization over \mathbf{y} is performed by the Viterbi algorithm.

(c) Jacob Eisenstein 2014-2017. Work in progress.

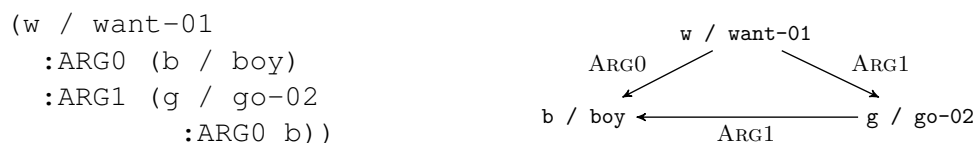


Figure 12.4: Two views of the AMR representation for the sentence *The boy wants to go*.

This model strongly outperformed alternative approaches at the time, including constrained decoding and convolutional neural networks. More recent work has combined recurrent neural network models with constrained decoding, using the *A** search algorithm to search over labelings that are feasible with respect to the constraints (He et al., 2017). This yields small improvements over the method of Zhou and Xu (2015). He et al. (2017) obtain larger improvements by creating an **ensemble** of SRL systems, each trained on an 80% subsample of the corpus. The average prediction across this ensemble is more robust than any individual model.

12.3 Abstract Meaning Representation

Semantic role labeling transforms the task of semantic parsing to a labeling task. Consider the sentence,

(12.22) The boy wants to go.

The PropBank semantic role labeling analysis is:

- (PREDICATE : *wants*, ARG0 : *the boy*, ARG1 : *to go*)
- (PREDICATE : *go*, ARG1 : *the boy*)

The **Abstract Meaning Representation (AMR)** unifies this analysis into a graph structure, in which each node is a **variable**, and each edge indicates a **concept** (Banarescu et al., 2013). This can be written in two ways, as shown in Figure 12.4. On the left is the PENMAN notation (Matthiessen and Bateman, 1991), in which each set of parentheses introduces a variable. Each variable is an **instance** of a concept, which is indicated with the slash notation: for example, *w / want-01* indicates that the variable *w* is an instance of the concept *want-01*, which in turn refers to the PropBank frame for the first sense of the verb *want*. Relations are introduced with colons: for example, *:arg0 (b / boy)* indicates a relation of type *arg0* with the newly-introduced variable *b*. Variables can be reused, so that when the variable *b* appears again as an argument to *g*, it is understood to refer to the same boy in both cases. This arrangement is indicated compactly in the graph structure on the right, with edges indicating concepts.

(c) Jacob Eisenstein 2014-2017. Work in progress.

AMR differs from PropBank-style semantic role labeling in a few key ways. First, it reifies each entity as a variable: for example, the *boy* in (12.22) is reified in the variable *b*, which is reused as ARG0 in its relationship with *w* / *want-01*, and as ARG1 in its relationship with *g* / *go-02*. Reifying entities as variables also makes it possible to represent the substructure of noun phrases more explicitly. For example, *Asha borrowed the algebra book* would be represented as:

```
(b / borrow-01
  :ARG0 (p / person
    :name (n / name
      :op1 "Asha"))
  :ARG1 (b2 / book
    :topic (a / algebra)))
```

This indicates that the variable *p* is a person, whose *name* is the variable *n*; that name has one token, the string *Asha*. Similarly, the variable *b2* is a book, and the *topic* of *b2* is a variable *a* whose type is *algebra*. The relations *name* and *topic* are examples of **non-core roles**, which are similar to adjunct modifiers in PropBank. However, AMR's inventory is more extensive, including more than 70 non-core roles, such as negation, time, manner, frequency, and location. Lists and sequences — such as the list of tokens in a name — are described using the roles *op1*, *op2*, etc.

Another key feature of AMR is that a semantic predicate can be introduced by any syntactic element. Consider the following examples, from Banarescu et al. (2013):

(12.23) The boy destroyed the room.

(12.24) the destruction of the room by the boy ...

(12.25) the boy's destruction of the room ...

All these examples have the same semantics in AMR,

```
(d / destroy-01
  :arg0 (b / boy)
  :arg1 (r / room))
```

Note that the noun *destruction* is linked to the verb *destroy*, which is captured by the PropBank frame *destroy-01*. This can happen with adjectives as well: in the phrase *the attractive spy*, the adjective *attractive* is linked to the PropBank frame *attract-01*:

```
(s / spy
  :arg0-of (a / attract-01))
```

(c) Jacob Eisenstein 2014-2017. Work in progress.

In this example, `arg0-of` is an **inverse relation**, indicating that `s` is the `arg0` of the predicate `a`. Inverse relations make it possible for all AMR parses to have a single root concept, which should be the **focus** of the utterance.

There are a number of other important linguistic issues in the design of AMR, which are summarized in the original paper (Banarescu et al., 2013) and the tutorial slides by Schneider et al. (2015). While AMR goes farther than semantic role labeling, it does not link semantically-related frames such as `buy/sell` (as FrameNet does), does not handle quantification (as first-order predicate calculus does), and makes no attempt to handle noun number and verb tense (as PropBank does). A recent survey by Abend and Rappoport (2017) situates AMR with respect to several other semantic representation schemes.

12.3.1 AMR Parsing

Abstract Meaning Representation is not a labeling of the original text — unlike PropBank semantic role labeling, and most of the other tagging and parsing tasks that we have encountered thus far. The AMR for a given sentence may include multiple concepts for single words in the sentence: as we have seen, the sentence *Asha likes algebra* contains both `person` and `name` concepts for the word *Asha*. Conversely, words in the sentence may not appear in the AMR: in *Boyang made a tour of campus*, the **light verb** *make* would not appear in the AMR, which would instead be rooted on the predicate `tour`. As a result, AMR is difficult to parse, and even evaluating AMR parsing involves considerable algorithmic complexity (Cai and Yates, 2013).

A further complexity is that AMR labeled datasets do not explicitly show the **alignment** between the AMR annotation and the words in the sentence. For example, the link between the word *wants* and the concept `want-01` is not annotated. To acquire training data for learning-based parsers, it is therefore necessary to first perform an alignment between the training sentences and their AMR parses. Flanigan et al. (2014) introduce a rule-based parser, which links text to concepts through a series of increasingly high-recall steps.

Graph-based parsing One family of approaches to AMR parsing is similar to the graph-based methods that we encountered in syntactic dependency parsing (chapter 10). For these systems (Flanigan et al., 2014), parsing is a two-step process:

1. **Concept identification** (Figure 12.5a). This involves constructing concept subgraphs for individual words or spans of adjacent words. For example, in the sentence, *Asha likes algebra*, we would hope to identify the minimal subtree including just the concept `like-01` for the word *like*, and the subtree `(p / person :name (n / name :op1 Asha))` for the word *Asha*.

(c) Jacob Eisenstein 2014-2017. Work in progress.

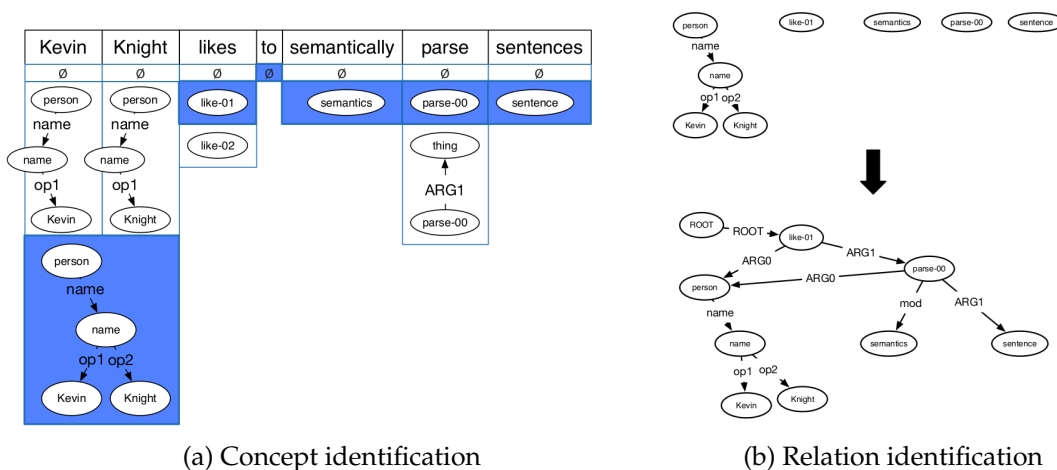


Figure 12.5: Subtasks for Abstract Meaning Representation parsing, from Schneider et al. (2015). [todo: ask for permission, or remake]

2. **Relation identification** (Figure 12.5b). This involves building a directed graph over the concepts, where the edges are labeled by the relation type. AMR imposes a number of constraints on the graph: all concepts must be included, the graph must be **connected** (there must be a path between every pair of nodes in the undirected version of the graph), and every node must have at most one outgoing edge of each type.

Both of these problems are solved by structure prediction. Concept identification requires simultaneously segmenting the text into spans, and labeling each span with a graph fragment containing one or more concepts. This is done by computing a set of features for each candidate span s and concept labeling c , and then returning the labeling with the highest overall score.

Relation identification can be formulated as search for the maximum spanning subgraph, under a set of constraints. Each labeled edge has a score, which is computed from features of the concepts. We then search for the set of labeled edges that maximizes the sum of these scores, under the constraint that the resulting graph is well-formed AMR. § 12.2.2 described how constrained optimization can be used for semantic role labeling; similar techniques have been applied to AMR relation identification (Flanigan et al., 2014).

Transition-based parsing In many cases, AMR parses are structurally similar to syntactic dependency parses. Figure 12.6 shows one such example. This motivates an alternative approach to AMR parsing: simply modify the syntactic dependency parse until it looks like a good AMR parse. Wang et al. (2015) propose a transition-based method, based on

(c) Jacob Eisenstein 2014-2017. Work in progress.

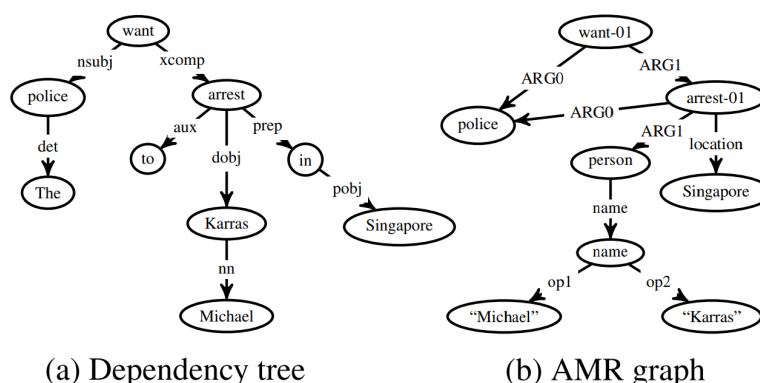


Figure 12.6: Syntactic dependency parse and AMR graph for the sentence *The police want to arrest Michael Karras in Singapore* (borrowed from Wang et al. (2015)) [todo: get permission]

incremental modifications to the syntactic dependency tree (you may review transition-based dependency parsing in § 10.3). At each step, the parser performs an action: for example, adding an AMR relation label to the current dependency edge, swapping the direction of a syntactic dependency edge, or cutting an edge and reattaching the orphaned subtree to a new parent. They train their system as a classifier, learning to choose the action as would be given by an **oracle** that is capable of reproducing the ground-truth parse. The 2016 SemEval evaluation compared a number of contemporary AMR parsing systems (May, 2016).

12.4 Applications of Predicate-Argument Semantics

Question answering **Factoid questions** have answers that are single words or phrases, such as *who discovered priors?*, *where was Barack Obama born?*, and *in what year did the Knicks last win the championship?* Shen and Lapata (2007) show that semantic role labeling can be used to answer such questions, by linking them to sentences in a corpus of text. They perform FrameNet semantic role labeling, making heavy use of dependency path features. For each sentence, they produce a weighted **bipartite graph**¹⁵ between FrameNet semantic roles and the words and phrases in the sentence. This is done by first scoring all pairs of semantic roles and assignments, as shown in the top half of Figure 12.8. They then find the bipartite edge cover, which is the minimum weighted subset of edges such that each vertex has at least one edge, as shown in the bottom half of Figure 12.8. After analyzing the question in this manner, Shen and Lapata then find semantically-compatible sentences in the corpus, by performing graph matching on the bipartite graphs for the question and

¹⁵A bipartite graph is one in which the vertices can be divided into two disjoint sets, and every edge connect a vertex in one set to a vertex in the other.

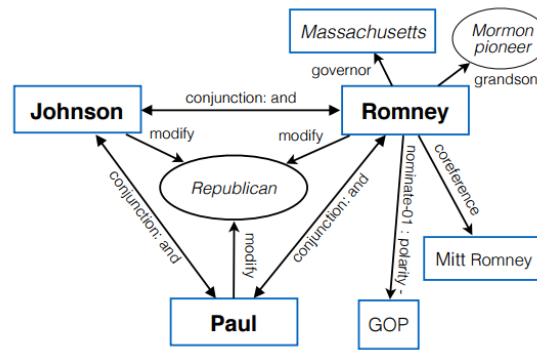


Figure 12.7: Fragment of AMR knowledge network for entity linking. Figure reprinted from Pan et al. (2015) [todo: permission]

candidate answer sentences. Finally, the **expected answer phrase** in the question — typically the *wh*-word — is linked to a phrase in the candidate answer source, and that phrase is returned as the answer.

Relation extraction The task of **relation extraction** involves identifying pairs of entities for which a given semantic relation holds. For example, we might like to find all $\langle i, j \rangle$ such that i is the **INVENTOR-OF** j . PropBank semantic role labeling can be applied to this task by identifying sentences whose verb signals the desired relation, and then extracting ARG1 and ARG2 as arguments. (To fully solve this task, these arguments must then be linked to entities, as described in chapter 16.) Christensen et al. (2010) compare the UIUC semantic role labeling system (which uses integer linear programming) against a simpler approach based on surface patterns (Banko et al., 2007). They find that the SRL system is considerably more accurate, but that it is several orders of magnitude slower. Conversely, Barnickel et al. (2009) apply SENNA, a convolutional neural network SRL system (Collobert and Weston, 2008) to the task of identifying biomedical relations (e.g., which genes inhibit or activate each other). In comparison with a strong baseline that applies a set of rules to syntactic dependency structures (Fundel et al., 2007), the SRL system is faster but less accurate. One possible explanation for these divergent results is that the Fundel et al. compare against a baseline which is carefully tuned for performance in a relatively narrow domain, while the system of Banko et al. is designed to analyze text across the entire web.

Entity linking Another core task in information extraction is to link mentions of entities (e.g., *Republican candidates like Romney, Paul, and Johnson ...*) to entities in a knowledge base (e.g., Lyndon Johnson or Gary Johnson). This is often done by examining nearby “collaborator” mentions — in this case, *Romney* and *Paul*. By jointly linking all

(c) Jacob Eisenstein 2014-2017. Work in progress.

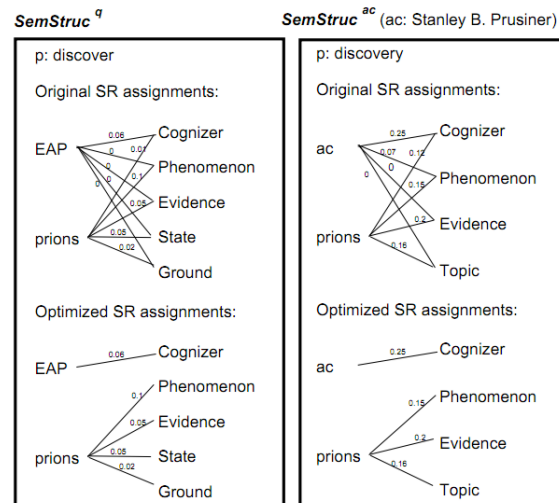


Figure 12.8: FrameNet semantic role labeling is used in factoid question answering, by aligning the semantic roles in the question (q) against those of sentences containing answer candidates (ac). “EAP” is the expected answer phrase, replacing the word *who* in the question. Figure reprinted from Shen and Lapata (2007) [todo: permission]

such mentions, it is possible to arrive at a good overall solution. Pan et al. (2015) apply AMR to this problem. For each entity, they construct a knowledge network based on its semantic relations with other mentions within the same sentence. They then rerank a set of candidate entities, based on the overlap between the entity’s knowledge network and the semantic relations present in the sentence (Figure 12.7). When applied to manually labeled AMR annotations, this approach is superior to state-of-the-art supervised methods that have access to labeled examples of linked mentions. Pan et al. also show that the method performs well from automated AMR, and that an AMR-based approach far outperforms a similar method based on PropBank semantic role labeling.[todo: rework for clarity]

Exercises

1. Write out an event semantic representation for the following sentences. You may make up your own predicates.

(12.26) *Abigail shares with Max.*

(12.27) *Abigail reluctantly shares a toy with Max.*

(12.28) *Abigail hates to share with Max.*

(c) Jacob Eisenstein 2014-2017. Work in progress.

2. Find the PropBank framesets for *share* and *hate* at <http://verbs.colorado.edu/propbank/framesets-english-aliases/>, and rewrite your answers from the previous question, using the thematic roles ARG0, ARG1, and ARG2.
3. Compute the syntactic path features for Abigail and Max in each of the example sentences (12.26) and (12.28) in Question 1, with respect to the verb *share*. If you're not sure about the parse, you can try an online parser such as <http://nlp.stanford.edu:8080/parser/>.
4. Compute the dependency path features for Abigail and Max in each of the example sentences (12.26) and (12.28) in Question 1, with respect to the verb *share*. Again, if you're not sure about the parse, you can try an online parser such as <http://nlp.stanford.edu:8080/parser/>. As a hint, the dependency relation between *share* and *Max* is OBL according to the Universal Dependency treebank (version 2).
5. PropBank semantic role labeling includes **reference arguments**, such as,

(12.29) $[_{AM-LOC}$ The bed] on $[_{R-AM-LOC}$ which] I slept broke.¹⁶

The label R-AM-LOC indicates that word *which* is a reference to *The bed*, which expresses the location of the event. Reference arguments must have referents: the tag R-AM-LOC can appear only when AM-LOC also appears in the sentence. Show how to express this as a linear constraint, specifically for the tag R-AM-LOC. Be sure to correctly handle the case in which neither AM-LOC nor R-AM-LOC appear in the sentence.

6. Explain how to express the constraints on semantic role labeling in Equation 12.12 and Equation 12.13 in the general form $\mathbf{A}z \geq b$.
7. Download the FrameNet sample data (<https://framenet.icsi.berkeley.edu/fndrupal/fulltextIndex>), and train a bag-of-words classifier to predict the frame that is evoked by each verb in each example. Your classifier should build a bag-of-words from the sentence in which the frame-evoking lexical unit appears. [todo: Somehow limit to one or a few lexical units.] [todo: use NLTK if possible]
8. Download the PropBank sample data, using NLTK (<http://www.nltk.org/howto/propbank.html>). Use a deep learning toolkit such as PyTorch or DyNet to train an LSTM to predict tags. You will have to convert the downloaded instances to a BIO sequence labeling representation first.
9. Produce the AMR annotations for the following examples:

¹⁶Example from 2013 NAACL tutorial slides by Shumin Wu

- (12.30) The girl likes the boy.
- (12.31) The girl was liked by the boy.
- (12.32) Abigail likes Maxwell Aristotle.
- (12.33) The spy likes the attractive boy.
- (12.34) The girl doesn't like the boy.
- (12.35) The girl likes her dog.

For (12.32), recall that multi-token names are created using `op1`, `op2`, etc. You will need to consult Banarescu et al. (2013) for (12.34), and Schneider et al. (2015) for (12.35). You may assume that *her* refers to *the girl* in this example.

10. Using an off-the-shelf PropBank SRL system,¹⁷ build a simplified question answering system in the style of Shen and Lapata (2007). Specifically, your system should do the following:
 - For each document in a collection, it should apply the semantic role labeler, and should store the output as a tuple.
 - For a question, your system should again apply the semantic role labeler. If any of the roles are filled by a *wh*-pronoun, you should mark that role as the expected answer phrase (EAP).
 - To answer the question, search for a stored tuple which matches the question as well as possible (same predicate, no incompatible semantic roles, and as many matching roles as possible). Align the EAP against its role filler in the stored tuple, and return this as the answer.

To evaluate your system, download a set of three news articles on the same topic, and write down five factoid questions that should be answerable from the articles. See if your system can answer these questions correctly. (If this problem is assigned to an entire class, you can build a large-scale test set and compare various approaches.)

¹⁷At the time of writing, the following systems are available: SENNA (<http://ronan.collobert.com/senna/>), Illinois Semantic Role Labeler (https://cogcomp.cs.illinois.edu/page/software_view/SRL), and mate-tools (<https://code.google.com/archive/p/mate-tools/>).

Chapter 13

Distributional and distributed semantics

A recurring theme in natural language processing is the complexity of the mapping from words to meaning. In chapter 3, we saw that a single word form, like *bank*, can have multiple meanings; conversely, a single meaning may be created by multiple surface forms, a lexical semantic relationship known as **synonymy**. Despite this complex mapping between words and meaning, natural language processing systems usually rely on words as the basic unit of analysis. This is especially true in semantics: the logical and frame semantic methods from the previous two chapters rely on hand-crafted lexicons that map from words to semantic predicates. But how can we analyze texts that contain words that we haven't seen before? This chapter describes methods that learn representations of word meaning by analyzing unlabeled data, vastly improving the generalizability of natural language processing systems. The theory that makes it possible to acquire meaningful representations from unlabeled data is the **distributional hypothesis**.

13.1 The distributional hypothesis

Here's a word you may not know: *tezgüino*.¹ If you do not know the meaning of *tezgüino*, then you are in the same situation as a natural language processing system when it encounters a word that did not appear in its training data. Now suppose you see that *tezgüino* is used in the following contexts:

- **C1:** *A bottle of _____ is on the table*
- **C2:** *Everybody likes _____.*
- **C3:** *Don't have _____ before you drive.*

¹The example is from Lin (1998).

	C1	C2	C3	C4	...
<i>tezguino</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	0	

Table 13.1: Distributional statistics for *tezguino* and five related terms

- **C4:** *We make _____ out of corn.*

What other words fit into these contexts? How about: *loud*, *motor oil*, *tortillas*, *choices*, *wine*? Each row of Table 13.1 is a vector that summarizes the contextual properties for each word, with a value of one for contexts in which the word can appear, and a value of zero for contexts in which it cannot. Based on these vectors, we can conclude:

- *wine* is very similar to *tezguino*;
- *motor oil* and *tortillas* are fairly similar to *tezguino*;
- *loud* is different.

These vectors, which we will call **word representations**, describe the **distributional** properties of each word. Does vector similarity imply semantic similarity? This is the **distributional hypothesis**, stated by Firth (1957) as: “You shall know a word by the company it keeps.” The distributional hypothesis has stood the test of time: distributional statistics are a core part of language technology today, because they make it possible to leverage large amounts of unlabeled data to learn about rare words that do not appear in labeled training data.

Distributional statistics have the striking ability to capture lexical semantic relationships such as analogies. Figure 13.1 shows two examples, based on two-dimensional projections of distributional **word embeddings**, discussed later in this chapter. In each case, word-pair relationships correspond to regular linear patterns in this two dimensional space. No labeled data about the nature of these relationships was required to identify this underlying structure.

Distributional semantics are computed from context statistics. **Distributed** semantics are a related but distinct idea: that meaning can be represented by numerical vectors rather than symbolic structures. Distributed representations are often estimated from distributional statistics, as in latent semantic analysis and WORD2VEC, described later in this

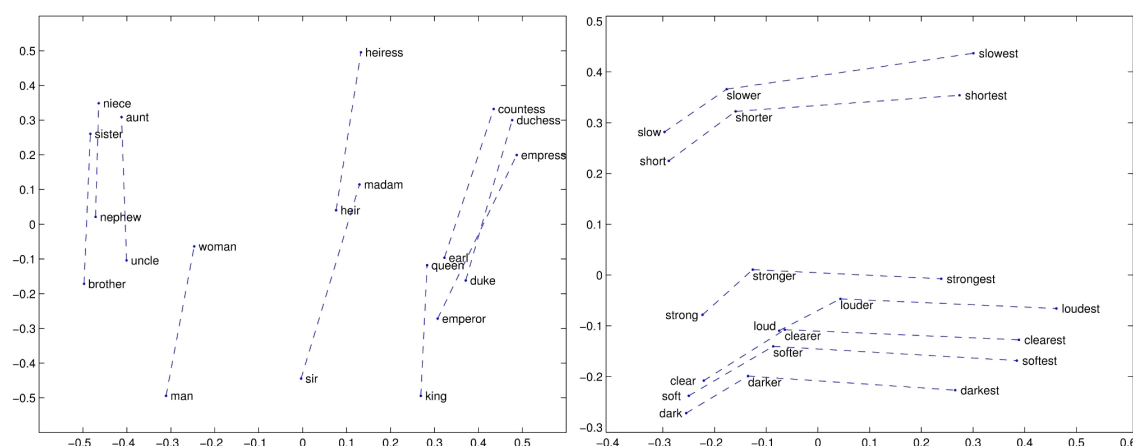


Figure 13.1: Lexical semantic relationships have regular linear structures in two dimensional projections of distributional statistics. From <http://nlp.stanford.edu/projects/glove/>. [todo: redo to make words bigger?]

chapter. However, distributed representations can also be learned in a supervised fashion from labeled data, as in the neural classification models encountered in chapter 2.

13.2 Design decisions for word representations

There are many approaches for computing word representations, but most can be distinguished on three main dimensions of variation: the nature of the representation, the source of contextual information, and the estimation procedure.

13.2.1 Representation

Today, the dominant word representations are k -dimensional vectors of real numbers, known as **word embeddings**. (The name is due to the fact that each discrete word is embedded in a continuous vector space.) This representation dates back at least to the late 1980s (Deerwester et al., 1990), and is used in popular techniques such as WORD2VEC (Mikolov et al., 2013a).

Dense vector word embeddings are well suited for neural network architectures, where they can be plugged in as inputs. They can also be applied in linear classifiers and structure prediction models (Turian et al., 2010), although it can be difficult to learn linear models that employ real-valued features (Kummerfeld et al., 2015). A popular alternative is bit-string representations, such as **Brown clusters** (§ 13.4), in which each word is represented by a variable-length sequence of zeros and ones (Brown et al., 1992).

(c) Jacob Eisenstein 2014-2017. Work in progress.

Another representational question is whether to estimate one embedding per surface form, or whether to estimate distinct embeddings for each word sense. Intuitively, if word representations are to capture the meaning of individual words, then words with multiple meanings should have multiple embeddings. This can be achieved by integrating unsupervised clustering with word embedding estimation (Huang and Yates, 2012; Li and Jurafsky, 2015). However, Arora et al. (2016) argue that it is unnecessary to model distinct word senses explicitly, because the embeddings for each surface form are a linear combination of the embeddings of the underlying senses.

13.2.2 Context

The distributional hypothesis says that word meaning is related to the “contexts” in which the word appears, but context can be defined in many ways. In the *tezgiino* example, contexts are entire sentences, but in practice there are far too many sentences for this to work. At the opposite extreme, the context could be defined as the immediately preceding word; this is the context considered in Brown clusters. WORD2VEC takes an intermediate approach, using local neighborhoods of words (e.g., $h = 5$) as contexts (Mikolov et al., 2013a). Contexts can also be much larger: for example, in **latent semantic analysis**, each word’s context vector includes an entry per document, with a value of one if the word appears in the document (Deerwester et al., 1990); in **explicit semantic analysis**, these documents are Wikipedia pages (Gabrilovich and Markovitch, 2007).

Words in context can be labeled by their position with respect to the target word w_m (e.g., two words before, one word after), which makes the resulting word representations more sensitive to syntactic differences (Ling et al., 2015a). Another way to incorporate syntax is to perform parsing as a preprocessing step, and then form context vectors from the set of dependency edges (Levy and Goldberg, 2014) or predicate-argument relations (Lin, 1998). The resulting context vectors for several of these methods are shown in Table 13.2.

The choice of context has a profound effect on the resulting representations, which can be viewed in terms of word similarity. Applying latent semantic analysis (§ 13.3) to contexts of size $h = 2$ and $h = 30$ yields the following nearest-neighbors for the word *dog*:²

- ($h = 2$): *cat, horse, fox, pet, rabbit, pig, animal, mongrel, sheep, pigeon*
- ($h = 30$): *kennel, puppy, pet, bitch, terrier, rottweiler, canine, cat, to bark, Alsatian*

²The example is from lecture slides by Marco Baroni, Alessandro Lenci, and Stefan Evert, who applied latent semantic analysis to the British National Corpus. You can find an online demo here: <http://clic.cimec.unitn.it/infomap-query/>

<i>The moment one learns English, complications set in.</i>	
Brown Clusters (Brown et al., 1992)	$\{one\}$
WORD2VEC (Mikolov et al., 2013a) ($h = 2$)	$\{moment, one, English, complications\}$
Structured WORD2VEC (Ling et al., 2015a) ($h = 2$)	$\{(moment, -2), (one, -1), (English, +1), (complications, +2)\}$
Dependency contexts (Levy and Goldberg, 2014)	$\{(one, NSUBJ), (English, DOBJ), (moment, ACL^{-1})\}$

Table 13.2: Contexts for the word *learns*, according to various word representations. For dependency context, $(one, NSUBJ)$ means that there is a relation of type NSUBJ (nominal subject) **to** the word *one*, and $(moment, ACL^{-1})$ means that there is a relation of type ACL (adjectival clause) **from** the word *moment*.

Which word list is better? Each word in the $h = 2$ list is an animal, reflecting the fact that locally, the word *dog* tends to appear in the same contexts as other animal types (e.g., *pet the dog*, *feed the dog*). In the $h = 30$ list, nearly everything is dog-related, including specific breeds such as *rottweiler* and *Alsatian*. The list also includes words that are not animals (*kennel*), and in one case (*to bark*), is not a noun at all. The 2-word context window is more sensitive to syntax, while the 30-word window is more sensitive to topic.

13.2.3 Estimation

Word embeddings are estimated by optimizing some objective: the likelihood of a set of unlabeled data (or a closely related quantity), or the reconstruction of a matrix of context counts, similar to Table 13.1.

Maximum likelihood estimation Likelihood-based optimization is derived from the objective $\log p(\mathbf{w}; \mathbf{u})$, where \mathbf{u}_i is the embedding of word i , and $\mathbf{w} = \{w_m\}_{m=1}^M$ is a corpus, represented as a list of M tokens. Recurrent neural network language models (§ 5.3) optimize this objective directly, backpropagating to the input word embeddings through the recurrent structure. However, state-of-the-art word embeddings employ huge corpora with hundreds of billions of tokens, and recurrent architectures are difficult to scale to such data. As a result, likelihood-based word embeddings are usually based on simplified likelihoods or heuristic approximations.

Matrix factorization The matrix $\mathbf{C} = \{\text{count}(i, j)\}$ stores the co-occurrence counts of word i and context j . Word representations can be obtained by approximately factoring this matrix, so that $\text{count}(i, j)$ is approximated by a function of a word embedding \mathbf{u}_i and

(c) Jacob Eisenstein 2014-2017. Work in progress.

a context embedding v_j . These embeddings can be obtained by minimizing the norm of the reconstruction error,

$$\min_{\mathbf{u}, \mathbf{v}} \|\mathbf{C} - \tilde{\mathbf{C}}(\mathbf{u}, \mathbf{v})\|_F, \quad [13.1]$$

where $\tilde{\mathbf{C}}(\mathbf{u}, \mathbf{v})$ is the approximate reconstruction resulting from the embeddings \mathbf{u} and \mathbf{v} , and $\|\mathbf{X}\|_F$ indicates the Frobenius norm, $\sum_{i,j} x_{i,j}^2$. Rather than factoring the matrix of word-context counts, it is useful to transform these counts using information-theoretic metrics such as **pointwise mutual information** (PMI), described in the next section.

13.3 Latent semantic analysis

Latent semantic analysis (LSA) is one of the oldest approaches to distributed semantics (Deerwester et al., 1990). It induces continuous vector representations of words by factoring a matrix of word and context counts, using **truncated singular value decomposition** (SVD),

$$\min_{\mathbf{U} \in \mathbb{R}^{V \times K}, \mathbf{S} \in \mathbb{R}^{K \times K}, \mathbf{V} \in \mathbb{R}^{|\mathcal{C}| \times K}} \|\mathbf{C} - \mathbf{U}\mathbf{S}\mathbf{V}^\top\|_F \quad [13.2]$$

$$s.t. \quad \mathbf{U}^\top \mathbf{U} = \mathbb{I} \quad [13.3]$$

$$\mathbf{V}^\top \mathbf{V} = \mathbb{I} \quad [13.4]$$

$$\forall i \neq j, \mathbf{S}_{i,j} = 0, \quad [13.5]$$

where V is the size of the vocabulary and $|\mathcal{C}|$ is the number of contexts. The word embeddings can then be set to the rows of the matrix \mathbf{U} . The matrix \mathbf{S} is constrained to be diagonal (these are called the singular values), and the columns of the product $\mathbf{S}\mathbf{V}^\top$ provide descriptions of the contexts. Each element $c_{i,j}$ is then reconstructed as a **bilinear product**,

$$c_{i,j} \approx \sum_{k=1}^K u_{i,k} s_k v_{j,k}, \quad [13.6]$$

and the objective is to minimize the sum of squared approximation errors. The orthonormality constraints $\mathbf{U}^\top \mathbf{U} = \mathbf{V}^\top \mathbf{V} = \mathbb{I}$ ensure that all pairs of dimensions in \mathbf{U} and \mathbf{V} are uncorrelated, so that each dimension conveys unique information. Efficient implementations of truncated singular value decomposition are available in numerical computing packages such as `scipy` and `matlab`.³

Latent semantic analysis is most effective when the count matrix is transformed before the application of SVD. One such transformation is **pointwise mutual information** (PMI;

³An important implementation detail is to represent \mathbf{C} as a **sparse matrix**, so that the storage cost is equal to the number of non-zero entries, rather than the size $V \times |\mathcal{C}|$.

Church and Hanks, 1990), which captures the degree of association between word i and context j ,

$$\text{PMI}(i, j) = \log \frac{p(i, j)}{p(i)p(j)} = \log \frac{p(i | j)p(j)}{p(i)p(j)} = \log \frac{p(i | j)}{p(i)} \quad [13.7]$$

$$= \log \text{count}(i, j) - \log \sum_{i'=1}^V \text{count}(i', j) - \log \sum_{j' \in \mathcal{C}} \text{count}(i, j') + \log \sum_{i'=1}^V \sum_{j' \in \mathcal{C}} \text{count}(i', j'), \quad [13.8]$$

The pointwise mutual information can be viewed as the logarithm of the ratio of the conditional probability of word i in context j to the marginal probability of word i in all contexts. When word i is statistically associated with context j , the ratio will be greater than one, so $\text{PMI}(i, j) > 0$. The PMI transformation focuses latent semantic analysis on reconstructing strong word-context associations, rather than on reconstructing large counts.

The PMI is negative when a word and context occur together less often than if they were independent, but such negative correlations are unreliable because counts of rare events have high variance. Furthermore, the PMI is undefined when $\text{count}(i, j) = 0$. **Positive PMI** (PPMI) is defined as,

$$\text{PPMI}(i, j) = \begin{cases} \text{PMI}(i, j), & p(i | j) > p(i) \\ 0, & \text{otherwise.} \end{cases} \quad [13.9]$$

Bullinaria and Levy (2007) compare a range of matrix transformations for latent semantic analysis, using a battery of tasks related to word meaning and word similarity (for more on evaluation, see § 13.6). They find that PPMI-based latent semantic analysis yields strong performance on a battery of tasks related to word meaning: for example, PPMI-based LSA vectors can be used to solve multiple-choice word similarity questions from the Test of English as a Foreign Language (TOEFL), obtaining 85% accuracy.

13.4 Brown clusters

Learning algorithms like perceptron and conditional random fields often perform better with discrete feature vectors. A simple way to obtain discrete representations from distributional statistics is by clustering (§ 4.1.1), so that words in the same cluster have similar distributional statistics. This can help in downstream tasks, by sharing features between all words in the same cluster. However, there is an obvious tradeoff: if the number of clusters is too small, the words in each cluster will not have much in common; if the number of clusters is too large, then the learner will not see enough examples from each cluster to generalize.

The solution to this problem is **hierarchical clustering**: using the distributional statistics to induce a tree-structured representation. Fragments of Brown cluster trees are shown

(c) Jacob Eisenstein 2014-2017. Work in progress.

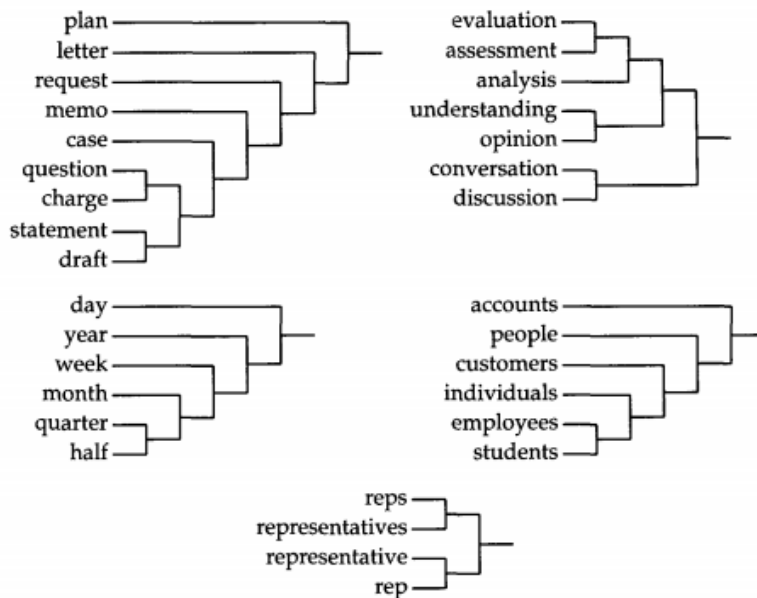


Figure 13.2: Some subtrees produced by bottom-up Brown clustering (Miller et al., 2004) on news text [todo: permission]

in Figure 13.2 and Table 13.3. Each word’s representation consists of a binary string describing a path through the tree: 0 for taking the left branch, and 1 for taking the right branch. In the subtree in the upper right of the figure, the representation of the word *conversation* is 01; the representation of the word *assessment* is 1110. Bitstring prefixes capture similarity at varying levels of specificity, and it is common to use the first eight, twelve, sixteen, and twenty bits as features in tasks such as named entity recognition (Miller et al., 2004) and dependency parsing (Koo et al., 2008).

Hierarchical trees can be induced from a likelihood-based objective, using a discrete latent variable $k_i \in \{1, 2, \dots, K\}$ to represent the cluster of word i :

$$\log p(\mathbf{w}; \mathbf{k}) \approx \sum_{m=1}^M \log p(w_m \mid w_{m-1}; \mathbf{k}) \quad [13.10]$$

$$\triangleq \sum_{m=1}^M \log p(w_m \mid k_{w_m}) + \log p(k_{w_m} \mid k_{w_{m-1}}). \quad [13.11]$$

This is similar to a hidden Markov model, with the crucial difference that each word can be emitted from only a single cluster: $\forall k \neq k_{w_m}, p(w_m \mid k) = 0$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

bitstring	ten most frequent words
01111010 0111	<i>excited thankful grateful stoked pumped anxious hyped psyched exited geeked</i>
01111010 100	<i>talking talkin complaining talkn bitching tlkn tlkin bragging rav- ing +k</i>
01111010 1010	<i>thinking thinkin dreaming worrying thinkn speakin reminiscing dreamin daydreaming fantasizing</i>
01111010 1011	<i>saying sayin suggesting stating sayn jokin talmbout implying insisting 5'2</i>
01111010 1100	<i>wonder dunno wondered duno donno dno dono wonda wounder dunnoe</i>
01111010 1101	<i>wondering wonders debating deciding pondering unsure won- derin debatin woundering wondern</i>
01111010 1110	<i>sure suree suuure suure sure- surre sures shuree</i>

Table 13.3: Fragment of a Brown clustering of Twitter data (Owoputi et al., 2013). Each row is a leaf in the tree, showing the ten most frequent words. This part of the tree emphasizes verbs of communicating and knowing, especially in the present participle. Each leaf node includes orthographic variants (*thinking*, *thinkin*, *thinkn*), semantically related terms (*excited*, *thankful*, *grateful*), and some outliers (5'2, +k). See http://www.cs.cmu.edu/~ark/TweetNLP/cluster_viewer.html for more.

Using the objective in Equation 13.11, the Brown clustering tree can be constructed from the bottom up: begin with each word in its own cluster, and incrementally merge clusters until only a single cluster remains. At each step, we merge the pair of clusters such that the objective in Equation 13.11 is maximized. Although the objective seems to involve a sum over the entire corpus, the score for each merger can be computed from the co-occurrence statistics $\text{count}(i, j)$, which counts the number of times a word from cluster i follows a word from cluster j . These counts can be updated incrementally as the clustering proceeds. The optimal merge at each step can be shown to maximize the **average mutual information**,

$$I(\mathbf{k}) = \sum_{k_1=1}^K \sum_{k_2=1}^K p(k_1, k_2) \times \text{PMI}(k_1, k_2) \quad [13.12]$$

$$p(k_1, k_2) = \frac{\text{count}(k_1, k_2)}{\sum_{k_1'=1}^K \sum_{k_2'=1}^K \text{count}(k_1', k_2')},$$

where $p(k_1, k_2)$ is the joint probability of a bigram involving a word in cluster k_1 followed by a word in k_2 . This probability and the PMI are both computed from the co-occurrence

Algorithm 15 Exchange clustering algorithm

```

procedure EXCHANGECLUSTERING( $\{\text{count}(\cdot, \cdot)\}, K$ )
  for  $i \in \mathcal{V}$  do                                     ▷ Compute unigram counts
     $\text{count}(i) \leftarrow \sum_{j \in \mathcal{C}} \text{count}(i, j)$ 
     $\{w_1, w_2, \dots, w_V\} \leftarrow \text{Sort-Descending}(\{\text{count}(\cdot)\})$ 
  for  $i \in 1 \dots K$  do                               ▷ Put each of the most common words in its own cluster
     $k_i \leftarrow i$ 
  for  $i \in (K + 1) \dots V$  do                         ▷ Iteratively add each word to the clustering
     $k_i \leftarrow K + 1$ 
    Let  $\langle k, k' \rangle$  be the two clusters whose merger maximizes  $I(\mathbf{k})$  (Equation 13.12)
    Renumber clusters so that  $k$  and  $k'$  are merged
    Update cluster co-occurrence counts
  repeat                                               ▷ Merge the remaining clusters into a tree
    Merge clusters  $\langle k, k' \rangle$  to maximize  $I(\mathbf{k})$ .
  until only one cluster remains
  return Tree structure induced by merge history

```

counts between clusters. After each merger, the co-occurrence vectors for the merged clusters are simply added up, so that the next optimal merger can be found efficiently.

This bottom-up procedure requires iterating over the entire vocabulary, and evaluating K_t^2 possible mergers at each step, where K_t is the current number of clusters at step t of the algorithm. Furthermore, computing the score for each merger involves a sum over K_t^2 clusters. The maximum number of clusters is $K_0 = V$, which occurs when every word is in its own cluster at the beginning of the algorithm. The time complexity is thus $\mathcal{O}(V^5)$.

To avoid this complexity, practical implementations use a heuristic approximation called **exchange clustering**. The K most common words are placed in clusters of their own at the beginning of the process. We then consider the next most common word, and merge it with one of the existing clusters. This continues until the entire vocabulary has been incorporated, at which point the K clusters are merged down to a single cluster, forming a tree. The algorithm never considers more than $K + 1$ clusters at any step, and the complexity is $\mathcal{O}(V \log V)$ in the size of the vocabulary, bounded by the cost of sorting the words at the beginning of the algorithm.

13.5 Neural word embeddings

Neural word embeddings combine aspects of the previous two methods: like latent semantic analysis, they are a continuous vector representation; like Brown clusters, they are trained from a likelihood-based objective. Let the vector \mathbf{u}_i represent the K -dimensional

(c) Jacob Eisenstein 2014-2017. Work in progress.

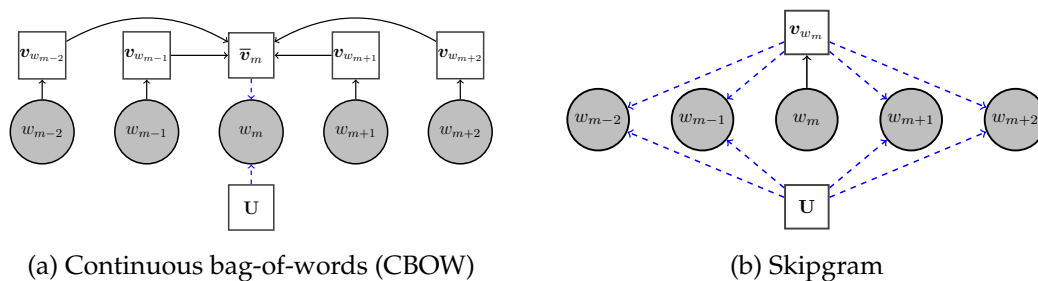


Figure 13.3: The CBOW and skipgram variants of WORD2VEC. The parameter \mathbf{U} is the matrix of word embeddings, and each \mathbf{v}_m is the context embedding for word w_m .

embedding for word i , and let \mathbf{v}_j represent the K -dimensional embedding for context j . The inner product $\mathbf{u}_i \cdot \mathbf{v}_j$ represents the compatibility between word i and context j . By incorporating this inner product into an approximation to the log-likelihood of a corpus, it is possible to estimate both parameters by backpropagation. WORD2VEC (Mikolov et al., 2013a) includes two such approximations: continuous bag-of-words (CBOW) and skipgrams.

13.5.1 Continuous bag-of-words (CBOW)

In recurrent neural network language models, each word w_m is conditioned on a recurrently-updated state vector, which is based on word representations going all the way back to the beginning of the text. The **continuous bag-of-words (CBOW)** model is a simplification: the local context is computed as an average of embeddings for words in the immediate neighborhood $m-h, m-h+1, \dots, m+h-1, m+h$,

$$\bar{\mathbf{v}}_m = \frac{1}{2h} \sum_{n=1}^h \mathbf{v}_{m+n} + \mathbf{v}_{m-n}. \quad [13.13]$$

Thus, CBOW is a bag-of-words model, because the order of the context words does not matter; it is continuous, because rather than conditioning on the words themselves, we condition on a continuous vector constructed from the word embeddings. The parameter h determines the neighborhood size, which Mikolov et al. (2013a) set to $h = 4$.

(c) Jacob Eisenstein 2014-2017. Work in progress.

The CBOW model optimizes an approximation to the corpus log-likelihood,

$$\log p(\mathbf{w}) \approx \sum_{m=1}^M \log p(w_m \mid w_{m-h}, w_{m-h+1}, \dots, w_{m+h-1}, w_{m+h}) \quad [13.14]$$

$$= \sum_{m=1}^M \log \frac{\exp(\mathbf{u}_{w_m} \cdot \bar{\mathbf{v}}_m)}{\sum_{j=1}^V \exp(\mathbf{u}_i \cdot \bar{\mathbf{v}}_m)} \quad [13.15]$$

$$= \sum_{m=1}^M \mathbf{u}_{w_m} \cdot \bar{\mathbf{v}}_m - \log \sum_{j=1}^V \exp(\mathbf{u}_i \cdot \bar{\mathbf{v}}_m). \quad [13.16]$$

13.5.2 Skipgrams

In the CBOW model, words are predicted from their context. In the **skipgram** model, the context is predicted from the word, yielding the objective:

$$\log p(\mathbf{w}) \approx \sum_{m=1}^M \sum_{n=1}^{h_m} \log p(w_{m-n} \mid w_m) + \log p(w_{m+n} \mid w_m) \quad [13.17]$$

$$= \sum_{m=1}^M \sum_{n=1}^{h_m} \log \frac{\exp(\mathbf{u}_{w_{m-n}} \cdot \mathbf{v}_{w_m})}{\sum_{j=1}^V \exp(\mathbf{u}_i \cdot \mathbf{v}_{w_m})} + \log \frac{\exp(\mathbf{u}_{w_{m+n}} \cdot \mathbf{v}_{w_m})}{\sum_{j=1}^V \exp(\mathbf{u}_i \cdot \mathbf{v}_{w_m})} \quad [13.18]$$

$$= \sum_{m=1}^M \sum_{n=1}^{h_m} \mathbf{u}_{w_{m-n}} \cdot \mathbf{v}_{w_m} + \mathbf{u}_{w_{m+n}} \cdot \mathbf{v}_{w_m} - 2 \log \sum_{j=1}^V \exp(\mathbf{u}_i \cdot \mathbf{v}_{w_m}). \quad [13.19]$$

In the skipgram approximation, each word is generated multiple times; each time it is conditioned only on a single word. This makes it possible to avoid averaging the word vectors, as in the CBOW model. The local neighborhood size h_m is randomly sampled from a uniform categorical distribution over the range $\{1, 2, \dots, h_{\max}\}$; Mikolov et al. (2013a) set $h_{\max} = 10$. Because the neighborhood grows outward with h , this approach has the effect of weighting near neighbors more than distant ones. Skipgram performs better on most evaluations than CBOW (see § 13.6 for details of how to evaluate word representations), but CBOW is faster to train Mikolov et al. (2013a).

13.5.3 Computational complexity

The WORD2VEC models can be viewed as an efficient alternative to recurrent neural network language models, which involves a recurrent state update whose time complexity is quadratic in the size of the recurrent state vector. CBOW and skipgram avoid this computation, and incur only a linear time complexity in the size of the word and context representations. However, all three models compute a normalized probability over

(c) Jacob Eisenstein 2014-2017. Work in progress.

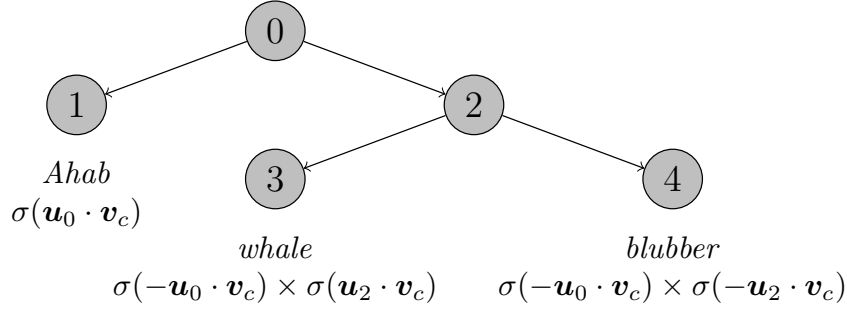


Figure 13.4: A fragment of a hierarchical softmax tree. The probability of each word is computed as a product of probabilities of local branching decisions in the tree.

word tokens; a naïve implementation of this probability requires summing over the entire vocabulary. The time complexity of this sum is $\mathcal{O}(V \times K)$, which dominates all other computational costs. There are two solutions: **hierarchical softmax**, a tree-based computation that reduces the cost to a logarithm of the size of the vocabulary; and **negative sampling**, an approximation that eliminates the dependence on vocabulary size. Both of these methods are also applicable to RNN language models.

Hierarchical softmax

In Brown clustering, the vocabulary is organized into a binary tree. Mnih and Hinton (2008) show that the normalized probability over words in the vocabulary can be reparametrized as a probability over paths through such a tree. This **hierarchical softmax** probability is computed as a product of binary decisions over whether to move left or right through the tree, with each binary decision represented as a sigmoid function of the inner product between the context embedding v_c and an output embedding associated with the node u_n ,

$$\Pr(\text{left at } n \mid c) = \sigma(u_n \cdot v_c) \quad [13.20]$$

$$\Pr(\text{right at } n \mid c) = 1 - \sigma(u_n \cdot v_c) = \sigma(-u_n \cdot v_c), \quad [13.21]$$

where σ refers to the sigmoid function, $\sigma(x) = \frac{1}{1 + \exp(-x)}$. The range of the sigmoid is the interval $(0, 1)$, and $1 - \sigma(x) = \sigma(-x)$.

As shown in Figure 13.4, the probability of generating each word is redefined as the product of the probabilities across its path. The sum of such probabilities is guaranteed to be one, for any context vector $v_c \in \mathbb{R}^K$. In a balanced binary tree, the depth is logarithmic in the number of leaf nodes, and thus the number of multiplications is equal to $\mathcal{O}(\log V)$. The number of non-leaf nodes is equal to $\mathcal{O}(2V - 1)$, so the number of parameters to be estimated increases by only a small multiple. The tree can be constructed using an incremental clustering procedure similar to hierarchical Brown clusters (Mnih and Hinton,

2008), or by using the Huffman (1952) encoding algorithm for lossless compression (Huffman, 1952).

Negative sampling

Likelihood-based methods are computationally intensive because each probability must be normalized over the vocabulary. These probabilities are based on scores for each word in each context, and it is possible to design an alternative objective that is based on these scores more directly: we seek word embeddings that maximize the difference between the score for the word that was really observed in each context, and the scores for a set of randomly selected **negative samples**:

$$\psi(i, j) = \log \sigma(\mathbf{u}_i \cdot \mathbf{v}_j) + \sum_{i' \in \mathcal{W}_{\text{neg}}} \log(1 - \sigma(\mathbf{u}_{i'} \cdot \mathbf{v}_j)), \quad [13.22]$$

where $\psi(i, j)$ is the score for word i in context j , σ refers to the sigmoid function, and \mathcal{W}_{neg} is the set of negative samples. The objective is to maximize the sum over the corpus, $\sum_{m=1}^M \psi(w_m, c_m)$, where w_m is token m and c_m is the associated context.

The set of negative samples \mathcal{W}_{neg} is obtained by sampling from a unigram language model. Mikolov et al. (2013b) construct this unigram language model by exponentiating the empirical word probabilities, setting $\hat{p}(i) \propto (\text{count}(i))^{\frac{3}{4}}$. This has the effect of redistributing probability mass from common to rare words. The number of negative samples increases the time complexity of training by a constant factor. Mikolov et al. (2013b) report that 5-20 negative samples works for small training sets, and that two to five samples suffice for larger corpora.

13.5.4 Word embeddings as matrix factorization

The negative sampling objective in Equation 13.22 can be justified as an efficient approximation to the log-likelihood, but it is also closely linked to the matrix factorization objective employed in latent semantic analysis. For a matrix of word-context pairs in which all counts are non-zero, negative sampling is equivalent to factorization of the matrix \mathbf{M} , where $M_{ij} = \text{PMI}(i, j) - \log k$: each cell in the matrix is equal to the pointwise mutual information of the word and context, shifted by $\log k$, with k equal to the number of negative samples (Levy and Goldberg, 2014). For word-context pairs that are not observed in the data, the pointwise mutual information is $-\infty$, but this can be addressed by considering only PMI values that are greater than $\log k$, resulting in a matrix of **shifted positive pointwise mutual information**,

$$M_{ij} = \max(0, \text{PMI}(i, j) - \log k). \quad [13.23]$$

Word embeddings are obtained by factoring this matrix with truncated singular value decomposition.

(c) Jacob Eisenstein 2014-2017. Work in progress.

word 1	word 2	similarity
<i>love</i>	<i>sex</i>	6.77
<i>stock</i>	<i>jaguar</i>	0.92
<i>money</i>	<i>cash</i>	9.15
<i>development</i>	<i>issue</i>	3.97
<i>lad</i>	<i>brother</i>	4.46

Table 13.4: Subset of the WS-353 (Finkelstein et al., 2002) dataset of word similarity ratings (examples from Faruqui et al. (2016b)).

GloVe (“global vectors”) are a closely related approach (Pennington et al., 2014), in which the matrix to be factored is constructed from log co-occurrence counts, $M_{ij} = \log \text{count}(i, j)$. We then minimize the sum of squares,

$$\begin{aligned} \min_{\mathbf{u}, \mathbf{v}, b, \tilde{b}} \quad & \sum_{j=1}^V \sum_{j \in \mathcal{C}} f(M_{ij}) \left(\widehat{\log M_{ij}} - \log M_{ij} \right)^2 \\ \text{s.t.} \quad & \widehat{\log M_{ij}} = \mathbf{u}_i \cdot \mathbf{v}_j + b_i + \tilde{b}_j, \end{aligned} \quad [13.24]$$

where b_i and \tilde{b}_j are offsets for word i and context j , which are estimated jointly with the embeddings \mathbf{u} and \mathbf{v} . The weighting function $f(M_{ij})$ is set to be zero at $M_{ij} = 0$, thus avoiding the problem of taking the logarithm of zero counts; it saturates at $M_{ij} = m_{\max}$, thus avoiding the problem of overcounting common word-context pairs. This heuristic turns out to be critical to the method’s performance.

The time complexity of sparse matrix reconstruction is determined by the number of non-zero word-context counts. Pennington et al. (2014) show that this number grows sublinearly with the size of the dataset: roughly $\mathcal{O}(N^{0.8})$ for typical English corpora. In contrast, the time complexity of WORD2VEC is linear in the corpus size. Computing the co-occurrence counts also requires linear time in the size of the corpus, but this operation can easily be parallelized using MapReduce-style algorithms (Dean and Ghemawat, 2008).

13.6 Evaluating word embeddings

Distributed word representations can be evaluated in two main ways. **Intrinsic** evaluations test whether the representations cohere with our intuitions about word meaning. **Extrinsic** evaluations test whether they are useful for downstream tasks, such as sequence labeling.

(c) Jacob Eisenstein 2014-2017. Work in progress.

13.6.1 Intrinsic evaluations

A basic question for word embeddings is whether the similarity of words i and j is reflected in the similarity of the vectors \mathbf{u}_i and \mathbf{u}_j . **Cosine similarity** is typically used to compare two word embeddings,

$$\cos(\mathbf{u}_i, \mathbf{u}_j) = \frac{\mathbf{u}_i \cdot \mathbf{u}_j}{\|\mathbf{u}_i\| \times \|\mathbf{u}_j\|}. \quad [13.25]$$

For any embedding method, we can then ask whether the cosine similarity of word embeddings is correlated with human judgments of word similarity. The WS-353 dataset (Finkelstein et al., 2002) includes similarity scores for 353 word pairs (Table 13.4). To test the accuracy of embeddings for rare and morphologically complex words, Luong et al. (2013) introduce a dataset of “rare words.” Outside of English, word similarity resources are limited, mainly consisting of translations of WS-353.

Word analogies (e.g., *king:queen :: man:woman*) have also been used to evaluate word embeddings (Mikolov et al., 2013c). In this evaluation, the system is provided with the first three parts of the analogy ($i_1 : j_1 :: i_2 : ?$), and the final element is predicted by finding the word embedding most similar to $\mathbf{u}_{i_1} - \mathbf{u}_{j_1} + \mathbf{u}_{i_2}$. Another evaluation tests whether word embeddings are related to broad lexical semantic categories called **supersenses** (Ciarrita and Johnson, 2003): verbs of motion, nouns that describe animals, nouns that describe body parts, and so on. These supersenses are annotated for English synsets in WordNet (Fellbaum, 2010). This evaluation is implemented in the `qvec` metric, which tests whether the matrix of supersenses can be reconstructed from the matrix of word embeddings (Tsvetkov et al., 2015).

Levy et al. (2015) compared several dense word representations for English — including latent semantic analysis, WORD2VEC, and GloVe — using six word similarity metrics and two analogy tasks. None of the embeddings outperformed the others on every task, but skipgrams were the most broadly competitive. Hyperparameter tuning played a key role: any method will perform badly if the wrong hyperparameters are used. Relevant hyperparameters include the embedding size, as well as algorithm-specific details such as the neighborhood size and the number of negative samples.

13.6.2 Extrinsic evaluations

Word representations contribute to downstream tasks like sequence labeling and document classification by enabling generalization across words. The use of distributed representations as features can be seen as a form of **semi-supervised learning**, in which performance on a supervised learning problem is augmented by learning distributed representations from unlabeled data (Miller et al., 2004; Koo et al., 2008; Turian et al., 2010). These **pre-trained word representations** can be used as features in a linear prediction model, or as the input layer in a neural network, such as a Bi-LSTM tagging model (§ 6.5.4). Word

(c) Jacob Eisenstein 2014-2017. Work in progress.

representations can be evaluated by the performance of the downstream systems that consume them: for example, GloVe embeddings are convincingly better than Latent Semantic Analysis as features in the downstream task of named entity recognition (Pennington et al., 2014). Unfortunately, extrinsic and intrinsic evaluations do not always point in the same direction, and the best word representations for one downstream task may perform poorly on another task (Schnabel et al., 2015).

When word representations are updated from labeled data in the downstream task, they are said to be **fine-tuned**. When labeled data is plentiful, pre-training may be unnecessary; when labeled data is scarce, fine-tuning may lead to overfitting. Various combinations of pre-training and fine-tuning can be employed. Pre-trained embeddings can be used as initialization before fine-tuning, and this can substantially improve performance (Lample et al., 2016). Alternatively, both fine-tuned and pre-trained embeddings can be used as inputs in a single model (Kim, 2014).

13.7 Distributed representations beyond distributional statistics

Distributional word representations can be estimated from huge unlabeled datasets, thereby covering many words that do not appear in labeled data: for example, GloVe embeddings are estimated from 800 billion tokens of web data,⁴ while the largest labeled datasets for NLP tasks are on the order of millions of tokens. Nonetheless, even a dataset of hundreds of billions of tokens will not cover every word that may be encountered in the future. Furthermore, many words will appear only a small number of times, making their embeddings unreliable. Many languages have more morphological structure than English, and thus have lower token-to-type ratios. When this problem is coupled with small training corpora, it becomes especially important to leverage other sources of information beyond distributional statistics.

13.7.1 Word-internal structure

One solution is to incorporate word-internal structure into word embeddings. Purely distributional approaches consider words as atomic units, but in fact, many words have internal structure, so that their meaning can be **composed** from the representations of sub-word units. Consider the following terms, all of which are missing from Google's pre-trained WORD2VEC embeddings:⁵

millicuries This word has **morphological** structure (see § 8.1.2 for more on morphology): the prefix *milli-* indicates an amount, and the suffix *-s* indicates a plural. (A *millicurie* is an unit of radioactivity.)

⁴<http://commoncrawl.org/>

⁵<https://code.google.com/archive/p/word2vec/>, accessed September 20, 2017

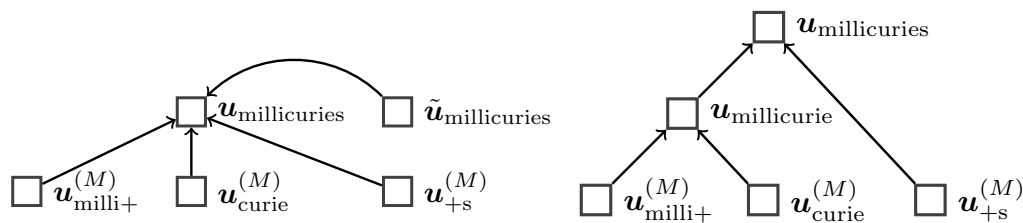


Figure 13.5: Two architectures for building word embeddings from subword units. On the left, morpheme embeddings $u^{(m)}$ are combined by addition with the non-compositional word embedding \tilde{u} (Botha and Blunsom, 2014). On the right, morpheme embeddings are combined in a recursive neural network (Luong et al., 2013).

caesium This word is a single morpheme, but the characters *-ium* are often associated with chemical elements. (*Caesium* is the British spelling of a chemical element, spelled *cesium* in American English.)

IAEA This term is an acronym, as suggested by the use of capitalization. The prefix *I-* frequently refers to international organizations, and the suffix *-A* often refers to agencies or associations. (IAEA is the International Atomic Energy Agency.)

Zhezghan This term is in title case, suggesting the name of a person or place, and the character bigram *zh* indicates that it is likely a transliteration. (*Zhezghan* is a mining facility in Kazakhstan.)

How can word-internal structure be incorporated into word representations? One approach is to construct word representations from embeddings of the characters or morphemes. For example, if word i has morphological segments \mathcal{M}_i , then its embedding can be constructed by addition (Botha and Blunsom, 2014),

$$u_i = \tilde{u}_i + \sum_{j \in \mathcal{M}_i} u_j^{(M)}, \quad [13.26]$$

where $u_m^{(M)}$ is a morpheme embedding and \tilde{u}_i is a non-compositional embedding of the whole word, which is an additional free parameter of the model (Figure 13.5, left side). All embeddings are estimated from a **log-bilinear language model** (Mnih and Hinton, 2007), which is similar to the CBOW model (§ 13.5), but includes only contextual information from preceding words. The morphological segments are obtained using an unsupervised segmenter (Creutz and Lagus, 2002). For words that do not appear in the training data, the embedding can be constructed directly from the morphemes, assuming that each morpheme appears in some other word in the training data. The free parameter \tilde{u} adds flexibility: words with similar morphemes are encouraged to have similar embeddings, but this parameter makes it possible for them to be different.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Subword information can be incorporated in various ways:

Subword units Examples like *IAEA* and *Zhezhgan* are not based on morphological composition, and a morphological segmenter is unlikely to identify meaningful subword units for these terms. Rather than using morphemes for subword embeddings, one can use characters (Santos and Zadrozny, 2014; Ling et al., 2015b; Kim et al., 2016), character n -grams (Wieting et al., 2016; Bojanowski et al., 2017), and **byte-pair encodings**, a compression technique which captures frequent substrings (Gage, 1994; Sennrich et al., 2016).

Composition Combining the subword embeddings by addition does not differentiate between orderings, nor does it identify any particular morpheme as the **root**. A range of more flexible compositional models have been considered, including recurrence (Ling et al., 2015b), convolution (Santos and Zadrozny, 2014; Kim et al., 2016), and **recursive neural networks** (Luong et al., 2013), in which representations of progressively larger units are constructed over a morphological parse, e.g. $((\text{milli}+\text{curie})+s)$, $((\text{in}+\text{flam})+\text{able})$, $(\text{in}+(\text{vis}+\text{ible}))$. A recursive embedding models is shown in the right panel of Figure 13.5.

Estimation Estimating subword embeddings from a full dataset is computationally expensive. An alternative approach is to train a subword model to match pre-trained word embeddings (Cotterell et al., 2016; Pinter et al., 2017). To train such a model, it is only necessary to iterate over the vocabulary, and not the corpus.

13.7.2 Lexical semantic resources

Resources such as WordNet provide another source of information about word meaning: if we know that *caesium* is a synonym of *cesium*, or that a *millicurie* is a type of *measurement unit*, then this should help to provide embeddings for the unknown words, and to smooth embeddings of rare words. One way to do this is to **retrofit** pre-trained word embeddings across a network of lexical semantic relationships, such as synonymy (Faruqui et al., 2015). This is done by minimizing the following objective,

$$\min_{\mathbf{U}} \sum_{j=1}^V \|\mathbf{u}_i - \hat{\mathbf{u}}_i\|_2 + \sum_{(i,j) \in \mathcal{L}} \beta_{ij} \|\mathbf{u}_i - \mathbf{u}_j\|_2, \quad [13.27]$$

where $\hat{\mathbf{u}}_i$ is the pretrained embedding of word i , and $\mathcal{L} = \{(i, j)\}$ is a lexicon of word relations. The hyperparameter β_{ij} controls the importance of adjacent words having similar embeddings; Faruqui et al. (2015) set it to the inverse of the degree of word i , $\beta_{ij} = |\{j : (i, j) \in \mathcal{L}\}|^{-1}$. Retrofitting improves performance on a range of intrinsic evaluations, and gives small improvements on an extrinsic document classification task.

(c) Jacob Eisenstein 2014-2017. Work in progress.

13.8 Distributed representations of multiword units

Can distributed semantics be extended to phrases, sentences, paragraphs, and beyond? Before exploring this possibility, recall the distinction between distributed and distributional representations. Neural embeddings such as WORD2VEC are both distributed (vector-based) and distributional (derived from counts of words in context). As we consider larger units of text, the counts in any corpus must decrease: in the limit, a multi-paragraph span of text would never appear twice, except in cases of plagiarism. Thus, the meaning of a large span of text cannot be determined from distributional statistics alone; it must be computed compositionally from smaller spans. But these considerations are orthogonal to the question of whether distributed representations — dense numerical vectors — are sufficiently expressive to capture the meaning of phrases, sentences, and paragraphs.

13.8.1 Purely distributional methods

Some multiword phrases are non-compositional: the meaning of such phrases is not derived from the meaning of the individual words using typical compositional semantics. This includes proper nouns like *San Francisco* as well as idiomatic expressions like *kick the bucket* (Baldwin and Kim, 2010). For these cases, purely distributional approaches can work. A simple approach is to identify multiword units that appear together frequently, and then treat these units as words, learning embeddings using a technique such as WORD2VEC. The problem of identifying multiword units is sometimes called **collocation extraction**, and can be approached using metrics such as pointwise mutual information: two-word units are extracted first, and then larger units are extracted. Mikolov et al. (2013b) identify such units and then treat them as words when estimating skipgram embeddings, showing that the resulting embeddings perform reasonably on a task of solving phrasal analogies, e.g. *New York : New York Times :: Baltimore : Baltimore Sun*.

13.8.2 Distributional-compositional hybrids

To move beyond short multiword phrases, composition is necessary. A simple but surprisingly powerful approach is to represent a sentence with the average of its word embeddings (Mitchell and Lapata, 2010). This can be considered a hybrid of the distributional and compositional approaches to semantics: the word embeddings are computed distributionally, and then the sentence representation is computed by composition. Baroni and Zamparelli (2010) embed short phrases by treating adjectives as matrices, which operate by multiplication on nouns.

The WORD2VEC approach can be stretched considerably further, embedding entire sentences using a model similar to skipgrams, in the “skip-thought” model of Kiros et al. (2015). Each sentence is **encoded** into a vector using a recurrent neural network: the encoding of sentence t is set to the RNN hidden state at its final token, $\mathbf{h}_{M_t}^{(t)}$. This vector is

(c) Jacob Eisenstein 2014-2017. Work in progress.

this was the only way
 it was the only way
 it was her turn to blink
 it was hard to tell
 it was time to move on
 he had to do it again
 they all looked at each other
 they all turned to look back
 they both turned to face him
they both turned and walked away

Figure 13.6: By interpolating between the distributed representations of two sentences (in bold), it is possible to generate grammatical sentences that combine aspects of both (Bowman et al., 2016)

then a parameter in a **decoder** model that is used to generate the previous and subsequent sentences: the decoder is another recurrent neural network, which takes the encoding of the neighboring sentence as an additional parameter in its recurrent update. (This **encoder-decoder model** is discussed at length in chapter 18.) The encoder and decoder are trained simultaneously from a likelihood-based objective, and the trained encoder can then be used to compute a distributed representation of any sentence. Skip-thought can be viewed as a hybrid of distributional and compositional approaches: the vector representation of each sentence is computed compositionally from the representations of the individual words, but the training objective is distributional, based on sentence co-occurrence across a corpus.

Autoencoders are a variant of encoder-decoder models in which the decoder is trained to produce the same text that was originally encoded, using only the distributed encoding vector (Li et al., 2015a). By interpolating between distributed representations of two sentences, $\alpha \mathbf{u}_i + (1 - \alpha) \mathbf{u}_j$, it is possible to generate sentences that combine aspects of the two inputs, as shown in Figure 13.6 (Bowman et al., 2016). Autoencoders can also be applied to longer texts, such as paragraphs and documents. This enables applications such as **question answering**, which can be performed by matching the encoding of the question with encodings of candidate answers (Miao et al., 2016).

13.8.3 Supervised compositional methods

Given a supervision signal, such as a label describing the sentiment or meaning of a sentence, a wide range of compositional methods can be applied to compute a distributed representation that can then be used to predict the label. The simplest is to average the embeddings of each word in the sentence, and then pass this average through a feed-

(c) Jacob Eisenstein 2014-2017. Work in progress.

forward neural network Iyyer et al. (2015). Convolutional and recurrent neural networks can effectively capture multiword phenomena such as negation (Kalchbrenner et al., 2014; Kim, 2014; Li et al., 2015b; Tang et al., 2015). Another approach is to incorporate the syntactic structure of the sentence. **Recursive neural networks** construct a representation for each syntactic constituent by operating on the representations of the children, thus propagating the distributed representation up the tree (Socher et al., 2012). However, Li et al. (2015b) show that in many cases, recurrent neural networks perform as well or better than recursive networks.

Whether convolutional, recurrent, or recursive, a key question is whether supervised sentence representations are task-specific, or whether a single supervised sentence representation model can yield useful performance on other tasks. Wieting et al. (2015) train a variety of sentence embedding models for the task of labeling pairs of sentences as **paraphrases**. They show that the resulting sentence embeddings give good performance for sentiment analysis. The **Stanford Natural Language Inference corpus** classifies sentence pairs as **entailments** (the truth of sentence i implies the truth of sentence j), **contradictions** (the truth of sentence i implies the falsity of sentence j), and **neutral** (i neither entails nor contradicts j). Sentence embeddings trained on this dataset transfer to a wide range of classification tasks (Conneau et al., 2017).

13.8.4 Hybrid distributed-symbolic representations

The power of distributed representations is in their generality: the distributed representation of a unit of text can serve as a summary of its meaning, and therefore as the input for downstream tasks such as classification, matching, and retrieval. For example, distributed sentence representations can be used to recognize the paraphrase relationship between closely related sentences like the following:

- (13.1) Donald thanked Vlad profusely.
- (13.2) Donald conveyed to Vlad his profound appreciation.
- (13.3) Vlad was showered with gratitude by Donald.

Symbolic representations are relatively brittle to this sort of variation, but are better suited to describe individual entities, the things that they do, and the things that are done to them. In examples (13.1)-(13.3), we not only know that somebody thanked someone else, but we can make a range of inferences about what has happened between the entities named *Donald* and *Vlad*. Because distributed representations do not treat entities symbolically, they lack the ability to reason about the roles played by entities across a sentence or larger discourse.⁶ A hybrid between distributed and symbolic representations

⁶At a 2014 workshop on semantic parsing, this critique of distributed representations was expressed by Ray Mooney — a leading researcher in computational semantics — in a now well-known quote, “you can’t cram the meaning of a whole sentence into a single vector!”

might give the best of both worlds: robustness to the many different ways of describing the same event, plus the expressiveness to support inferences about entities and the roles that they play.

A “top-down” hybrid approach is to begin with logical semantics (of the sort described in the previous two chapters), and then relax the requirement that a predefined lexicon associate words with each of a large set of predefined semantic predicates: instead, predicates are identified with distributional word clusters (Poon and Domingos, 2009; Lewis and Steedman, 2013). A “bottom-up” approach is to add minimal symbolic structure to existing distributed representations, such as vector representations for each entity (Ji and Eisenstein, 2015; Wiseman et al., 2016). This has been shown to improve performance on two problems that we will encounter in the following chapters: classification of **discourse relations** between adjacent sentences (chapter 15; Ji and Eisenstein, 2015), and **coreference resolution** of entity mentions (chapter 14; Wiseman et al., 2016; Ji et al., 2017). Research on hybrid semantic representations is still in an early stage, and future representations may deviate more boldly from existing symbolic and distributional approaches.

Additional reading

Turney and Pantel (2010) survey a number of facets of vector word representations, focusing on matrix factorization methods. Schnabel et al. (2015) highlight problems with similarity-based evaluations of word embeddings, and present a novel evaluation that controls for word frequency. Baroni et al. (2014) address linguistic issues that arise in attempts to combine distributed and compositional representations.

In bilingual and multilingual distributed representations, embeddings are estimated for translation pairs or tuples, such as $\langle \text{dog}, \text{perro}, \text{chien} \rangle$. These embeddings can improve machine translation (Zou et al., 2013; Klementiev et al., 2012), transfer natural language processing models across languages (Täckström et al., 2012), and make monolingual word embeddings more accurate (Faruqui and Dyer, 2014). A typical approach is to learn a projection that maximizes the correlation of the distributed representations of each element in a translation pair, which can be obtained from a bilingual dictionary. Distributed representations can also be linked to perceptual information, such as image features. Bruni et al. (2014) use textual descriptions of images to obtain visual contextual information for various words, which supplements traditional distributional context. Image features can also be inserted as contextual information in log bilinear language models (Kiros et al., 2014), making it possible to automatically generate text descriptions of images.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Exercises

1. Prove that the sum of probabilities of paths through a hierarchical softmax tree is equal to one.
2. In skipgram word embeddings, the negative sampling objective can be written as,

$$\mathcal{L} = \sum_{i \in \mathcal{V}} \sum_{j \in \mathcal{C}} \text{count}(i, j) \psi(i, j), \quad [13.28]$$

with $\psi(i, j)$ is defined in Equation 13.22.

Suppose we draw the negative samples from the empirical unigram distribution $\hat{p}(i) = p_{\text{unigram}}(i)$. First, compute the expectation of \mathcal{L} with respect this probability.

Next, take the derivative of this expectation with respect to the score of a single word context pair $\sigma(\mathbf{u}_i \cdot \mathbf{v}_j)$, and solve for the pointwise mutual information $\text{PMI}(i, j)$. You should be able to show that at the optimum, the PMI is a simple function of $\sigma(\mathbf{u}_i \cdot \mathbf{v}_j)$ and the number of negative samples.

3. * In Brown clustering, prove that the cluster merge that maximizes the average mutual information (Equation 13.12) also maximizes the log-likelihood objective (Equation 13.11).

For the next two problems, download a set of pre-trained word embeddings, such as the WORD2VEC or polyglot embeddings.

4. Use cosine similarity to find the most similar words to: *dog*, *whale*, *before*, *however*, *fabricate*.
5. Use vector addition and subtraction to compute target vectors for the analogies below. After computing each target vector, find the top three candidates by cosine similarity.
 - *dog:puppy :: cat: ?*
 - *speak:speaker :: sing:?*
 - *France:French :: England:?*
 - *France:wine :: England:?*

The remaining problems will require you to build a classifier and test its properties. Pick a multi-class text classification dataset, such as RCV1⁷). Divide your data into training (60%), development (20%), and test sets (20%), if no such division already exists.

⁷http://www.ai.mit.edu/projects/jmlr/papers/volume5/lewis04a/lyrl2004_rcv1v2_README.htm

6. Train a convolutional neural network, with inputs set to pre-trained word embeddings from the previous problem. Use a special, fine-tuned embedding for out-of-vocabulary words. Train until performance on the development set does not improve. You can also use the development set to tune the model architecture, such as the convolution width and depth. Report F -measure and accuracy, as well as training time.
7. Now modify your model from the previous problem to fine-tune the word embeddings. Report F -measure, accuracy, and training time.
8. Try a simpler approach, in which word embeddings in the document are averaged, and then this average is passed through a feed-forward neural network. Again, use the development data to tune the model architecture. How close is the accuracy to the convolutional networks from the previous problems?

Chapter 14

Reference Resolution

References are one of the most noticeable forms of linguistic ambiguity, afflicting not just automated natural language processing systems, but also fluent human readers. For this reason, warnings to avoid “ambiguous pronouns” are ubiquitous in manuals and tutorials on writing style. But referential ambiguity is not limited to pronouns, as shown in the following text:

(14.1) *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... **[Cook]**₄ is on **[his]**₅ first trip to **[the country]**₆ since taking over...*¹

Each of the bolded substrings in the passage refers to an entity that is introduced earlier in the story. These references include the pronouns *he* and *his*, but also the shortened name *Cook*, and most challengingly, **nominals** such as *the firm* and *the firm’s biggest growth market*. Only by resolving several of these references can we reach the (correct) inference that China is Apple’s biggest growth market.

The task of reference resolution is often broken into two components:

- **Coreference resolution**, which is the task of linking spans of text such as the *the firm* to other spans, such as *Apple Inc*. A subset of coreference resolution is the task of **anaphora resolution**, which typically involves resolving only pronoun anaphora such as *he* and *her*.
- **Entity linking**,² which is the task of linking spans of text to entities in a knowledge base. This step is a prerequisite for the model-based semantic parsing that was

¹<http://www.reuters.com/article/us-apple-china-idUSBRE82Q06420120327>, retrieved on March 26, 2017

²Amusingly, there are many names for this task: deduplication, approximate string match, entity clustering, record linking, multidocument coreference resolution, etc.

considered in chapter 11.

These tasks have traditionally been distinguished because they seem to require different sorts of knowledge to perform, and different resources to evaluate. As we will see, coreference resolution — especially anaphora resolution — is constrained by syntax and by compatibility of attributes such as gender, number, and animacy. Coreference resolution can be evaluated by comparing against a ground truth that is specified at the document level, without reference to any external information. In contrast, solving entity linking requires making inferences about name compatibility and about semantic properties of each entity — although these inferences are sometimes necessary for coreference resolution too. Evaluation of entity linking requires linking textual references to some predefined external ontology. Of the two tasks, research on coreference resolution is more mature, and will therefore be the focus of this chapter. Entity linking is also considered as a subtask of information extraction, and is described in § 16.1.

14.1 Forms of referring expressions

The three main forms of referring expressions — pronouns, names, and nominals — each pose unique challenges for the reader. As a coarse-grained summary, pronouns are constrained by syntax and semantic attributes; name references constrained by rules for matching; nominals are linked by world knowledge.

14.1.1 Pronouns

Pronouns are a closed class of words that are used for references. A natural way to think about pronoun resolution is what Kehler (2007) calls the SMASH approach:

- **S**earch for candidate referents;
- **M**atch against hard agreement constraints;
- **A**nd **S**elect using **H**euristics, which are “soft” constraints such as recency and parallelism.

In the search step, candidates are identified from the preceding text or speech.³ In models such as **centering theory**, any entity that has previously been **evoked** can be **accessed** in any subsequent unit of text (Grosz et al., 1995). However, cognitive constraints may imply

³Pronouns whose referents come later are known as **cataphora**, e.g.,

(14.1) *Many years later, as **he** faced the firing squad, **Colonel Aureliano Buendía** was to remember that distant afternoon when his father took him to discover ice.*

(This is the first sentence of *One Hundred Years of Solitude*, by Gabriel García Márquez.)

(c) Jacob Eisenstein 2014-2017. Work in progress.

that entities which have not been mentioned recently are unlikely to be accessed without be re-introduced, and correspondingly, computational constraints may encourage algorithms to consider only referents that are relatively recent.

Matching constraints for pronouns

Semantic constraints include morphologically marked information such as number, person, gender, and animacy.

- (14.2) Tim Cook has jetted in for talks with officials as **he** seeks to clear up a pile of problems...

We can identify the following features of the pronoun and possible referents:

- Number(*he*) = singular
- Number(*officials*) = plural
- Number(*Tim Cook*) = singular

Since there are no other possible referents, *he* almost certainly refers back to *Tim Cook*.

Other features include person, gender, and animacy, as in the following examples:

- (14.3) *Sally met my brother. He charmed her.*

- (14.4) *Sally met my brother. She charmed him.*

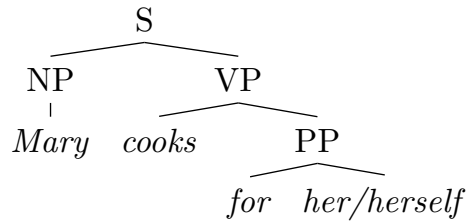
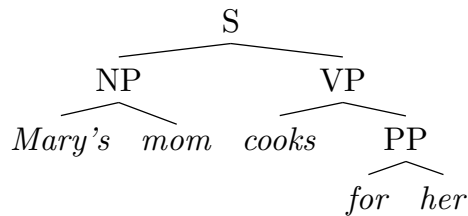
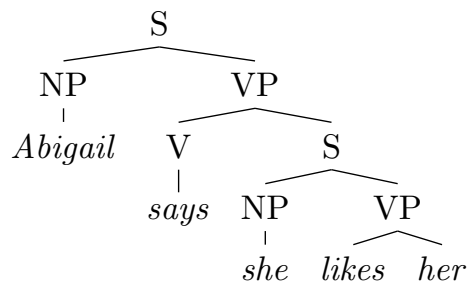
- (14.5) **We₁ told them₁ not to go.*

- (14.6) *Putin brought a bottle of vodka. It was from Russia.*

Another source of constraints comes from syntax. To understand these constraints, it is helpful to introduce some linguistic terminology:

- 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**.

For example, consider the tree in Figure 14.1. In this example, *Mary* c-commands *her/herself*, because the first branching node above *Mary* also dominates *her/herself*. However, *her/herself* does not c-command *Mary*. Thus, the pronoun *her* **cannot** refer to *Mary*, because non-reflexive pronouns cannot refer to antecedents that c-command them. On the other hand, the reflexive *herself* **must** refer to *Mary*.

Figure 14.1: *Mary* c-commands *her/herself*Figure 14.2: *Mary* does not c-command *her*, but *Mary's mom* does.Figure 14.3: The scope of *Abigail* is limited by the S non-terminal. Either *she* or *her* (but not both) can bind to *Abigail*.

Now consider the example, shown in Figure 14.2. Here, *Mary* does **not** c-command *her*, but *Mary's mom* c-commands *her*. Thus, *her* **can** refer to *Mary* — and we cannot use reflexive *herself* in this context, unless we are talking about *Mary's mom*. However, *her* does not have to refer to *Mary*.

A more complex example is shown in Figure 14.3. This indicates how the constraints defined here have a limited domain. [**todo: explain how this is limited**] The pronoun *she* can refer to *Abigail*, because *Abigail* is outside the domain of *she*. Similarly, *her* can also refer to *Abigail*. But *she* and *her* cannot be coreferent.

Heuristics

After applying constraints, there will be a number of candidate referents for each pronoun. In the SMASH paradigm, heuristics are then applied to compare among these possibilities.

Recency is a particularly strong heuristic. All things equal, readers will prefer the more recent referent for a given pronoun, particularly when comparing referents that occur in different sentences. Jurafsky and Martin (2009) offer the following example:

(14.7) *The doctor found an old map in the captain's chest. Jim found an even older map hidden on the shelf. It described an island.*

Readers are expected to prefer the second, older map as the referent for the pronoun *it*.

However, subjects are often preferred over objects, and this can contradict the preference for recency when two candidate referents are in the same sentence. For example,

(14.8) *Asha loaned Mei a book on Spanish. She is always trying to help people.*

Here, we may prefer to link *she* to *Asha* rather than *Mei*, because of *Asha*'s position in the subject role of the preceding sentence. (Arguably, this preference would be reversed if the second sentence were *She is visiting Argentina next month.*)

A third heuristic is parallelism:

(14.9) *Asha loaned Mei a book on Spanish. Olya loaned her a book on Portuguese.*

Here *Mei* is preferred as the referent for *her*, contradicting the preference for the subject *Asha* in the preceding sentence.

Hobbs (1978) unifies recency and subject-role heuristics by traversing the document in a syntax-driven fashion: each preceding sentence is traversed breadth-first, left-to-right (Figure 14.4). In this way, *Asha* would be preferred as the referent for *she* in (14.8), while the older map would be preferred as the referent for *it* in *ex:coref-recency*. **Centering theory** offers an alternative unification of recency and syntactic prominence, maintaining ordered lists of candidate referents throughout the text or discourse (Grosz et al., 1995). Centering can also be viewed as a generative model, in that it predicts the form of the referring expression that will be used for each entity reference in a sentence.

In early work on reference resolution, Lappin and Leass (1994) set weights for a half-dozen syntactic preferences by hand, choosing the referent with the highest overall weight. More recent work uses machine learning approaches to quantify the importance of each of these factors, as discussed in § 14.2. However, pronoun resolution often cannot be performed successfully using syntactic heuristics alone. This is shown by the classic example pair:

(c) Jacob Eisenstein 2014-2017. Work in progress.

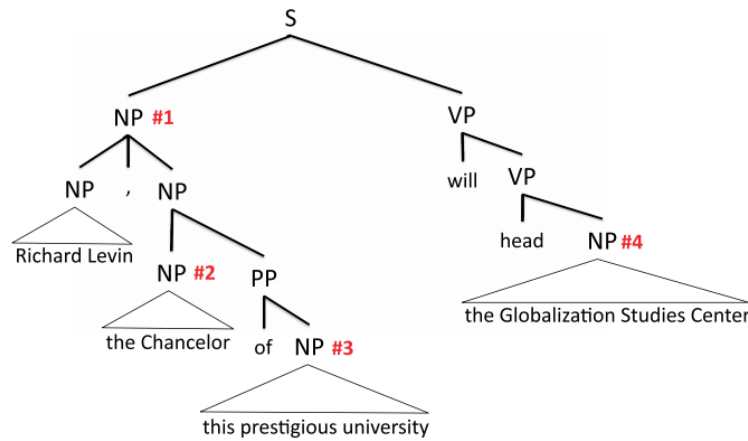


Figure 14.4: Left-to-right breadth-first tree traversal, proposed by Hobbs (1978), as implemented by Lee et al. (2013)

(14.10) *The **city council** denied the protesters a permit because **they** feared violence.*

(14.11) *The city council denied **the protesters** a permit because **they** advocated violence.*⁴

Non-referential pronouns

While pronouns are generally used to refer to things, they need not refer to entities, as shown in the following examples:

(14.12) *They told me that I was too ugly, but I didn't believe **it**.*

(14.13) *Alice saw Bob get angry, and I saw **it** too.*

(14.14) *They told me that I was too ugly, but **that** was a lie.*

(14.15) *Jess said she worked in security.
I suppose **that's** one way to put it.*

Pronouns may also have **generic** referents, meaning that they do not refer to entities in any model, but rather, to possible entities:

(14.16) *A good father takes care of **his** kids.*

(14.17) *On the moon, **you** have to carry **your** own oxygen.*

(14.18) *Every farmer who owns a donkey beats it. (Geach, 1962)*

⁴This pair is attributed to Winograd (1972), but I downloaded that article and didn't find it.

Finally, pronouns need not refer to anything at all:

- (14.19) *It's raining.*
- (14.20) *It's crazy out there.*
- (14.21) *It's money that she's really after.*
- (14.22) *It sucks that we have to work so hard.*

In the first two examples above, *it* is **pleonastic**; the third and fourth examples are **cleft** and **extraposition**. How can we automatically distinguish these usages of *it* from referential pronouns? Bergsma et al. (2008) propose a substitutability test. Consider the difference between the following two examples:

- (14.23) *You can make it in advance.*
- (14.24) *You can make it in showbiz.*

In the second example, the pronoun **it** is non-referential. One way to see this is by substituting another pronoun, like **them**, into these examples:

- (14.25) *You can make them in advance.*
- (14.26) *?You can make them in showbiz.*

The questionable grammaticality of the second example suggests that **it** cannot be referential. Bergsma et al. (2008) operationalize this idea by comparing distributional statistics for 5-grams around the word *it*, testing how often other pronouns or nouns ever appear in the same position as *it*. In cases where other pronouns are frequent, the *it* is likely referential.

14.1.2 Proper Names

If a proper name is used as a referring expression, it often refers to another proper name, so that the coreference problem is simply to determine whether the two names match. Subsequent proper name references often use a shortened form, as in the running example:

- (14.27) *Apple Inc Chief Executive **Tim Cook** has jetted into China ... **Cook** is on his first business trip to the country since taking over ...*

In this news article, the family name *Cook* is used as a referring expression; in informal conversation, it might be more typical to use the given name *Tim*, while more formal venues, such as *The Economist*, would use the title form *Mr Cook*.

Thus, exact match is unlikely to identify many proper name references. A typical solution is to match the syntactic **head words** of the reference with the referent. Recall that the head word of a phrase can be identified by applying head percolation rules to the phrasal parse tree (§ 9.4.2); alternatively, the head can be identified as the root of the dependency subtree covering the name. For sequences of proper nouns, the head word will be the final word, which in the example is *Cook*.

While useful, there are a number of caveats to the practice of matching head words of proper names.

- In the European tradition, family names tend to be more specific than given names, and family names usually come last. However, other traditions have other practices: for example, in Chinese names, the family name typically comes first; in Japanese, honorifics come after the name, as in *Nobu-San* (*Mr. Nobu*).
- In many organization names, it is also the case that the head word is not the most informative: for example, *Georgia Tech* and *Virginia Tech* are distinguished by the modifiers *Virginia* and *Georgia*, and not the heads. This concern applies even when the referring expression is a substring of the candidate referent: *Lebanon* does not refer to the same entity as *Southern Lebanon*, and Lee et al. (2011) add a rule to deal with the specific case of geographical modifiers.
- Finally, proper names can be nested, as in [*the CEO of [Microsoft]*]. Haghighi and Klein (2009) introduce a constraint to prevent nested noun phrases from being marked as coreferential.

Despite these difficulties, proper names are the easiest category of references to resolve (Stoyanov et al., 2009). In machine learning systems, one solution is to include a range of matching features, including exact match, head match, string inclusion, and even matching on “bags” of tokens, so that, e.g., *Tim Cook* matches *Cook*, *Tim* (Bontcheva et al., 2002)[**todo: can't find this citation in TALN proceedings**]. In addition to matching features, competitive systems include large lists, or **gazetteers**, of acronyms (e.g., *the National Basketball Association/NBA*), demonyms (e.g., *the Israelis/Israel*), and other aliases (e.g., *the Georgia Institute of Technology/Georgia Tech*).⁵ The learning algorithm can then determine the appropriate weights for each matching feature.

14.1.3 Nominals

In coreference resolution, noun phrases that are neither pronouns nor names are referred to as **nominals**. In the running example, nominal references include:

⁵Lists of aliases were used heavily in the Message Understanding Conference (MUC) systems of the 1990s, which helped to define the coreference resolution task (Grishman and Sundheim, 1996). They are still used in some of the most competitive systems at the time of this writing (e.g. Martschat and Strube, 2015).

- *the firm (Apple Inc)*
- *the firm's biggest growth market (China)*
- *the country (China).*

Nominals are generally more difficult to resolve than pronouns and names (Durrett and Klein, 2013, e.g.), and the examples above suggest why this may be the case: world knowledge is required to identify *Apple Inc* as a *firm*, and *China* as a *growth market*. Other difficult examples include the use of colloquial expressions, such as coreference between *Clinton transition officials* and *the Clinton camp* (Soon et al., 2001). But there are also cases that can be handled by surface features such as head word match: for example, *the tax cut bill* may be referenced later by *the Republican bill* or even *the bill*.

Attempts to use semantics to improve nominal coreference have met with limited success. Durrett and Klein (2013) employ WordNet synonymy and hypernymy relations on head words, named entity types (e.g., person, organization), and unsupervised clustering over nominal heads. These features give only limited improvement over simple baseline using surface features such as string match.

14.2 Learning for coreference resolution

Coreference resolution is a non-traditional learning problem, because it is not obvious what constitutes an “instance.” A number of proposals have been put forward. In discussing these approaches, references and candidate referents are called **mentions**; chains of references are called **entities**.⁶ [todo: add figure] Ground truth annotations identify the entities. In our running example, this would be:

- *Apple Inc Chief Executive Tim Cook, he, Cook*
- *Apple Inc, the firm*
- *China, the firm's biggest growth market, the country.*

“Singleton” mentions (e.g., *government officials*) are annotated in the ACE Corpus⁷, but not in the OntoNotes corpus (Hovy et al., 2006).

Coreference resolution can be viewed as a structure prediction problem, where the goal is to identify a set of coreference chains c among all possible coreference structures

⁶In many annotations, the term **markable** is used to refer to spans of text that can **potentially** mention an entity. The set of markables includes non-referential pronouns such as pleonastic *it*, which does not mention any entity. Part of the job of the coreference system is to avoid incorrectly linking these non-referential markables to any mention chains.

⁷<https://www.ldc.upenn.edu/collaborations/past-projects/ace/annotation-tasks-and-specifications>

$\mathcal{C}(w)$,

$$\operatorname{argmax}_{c \in \mathcal{C}(w)} \theta \cdot f(c, w). \quad [14.1]$$

Each chain c_i consists of a set of mentions $\{m_j\}$. Typically it is the coreference resolution system's job to identify the mentions from unannotated text, although systems are sometimes evaluated with "gold" mentions from the annotators.

The main approaches to coreference resolution can be distinguished by how they decompose the feature function $f(c, w)$. In mention-based models, features are defined over pairs of entities. This can facilitate inference, but mention-based models can suffer from incoherent entity chains, such as $\{Hillary Clinton \leftarrow Clinton \leftarrow Mr Clinton\}$. In entity-based models, the goal is to ensure that the entire entity chain is coherent. This can make inference more difficult, since the number of possible entity groupings is exponential in the number of mentions. A second distinction is whether the training instances are pairs (mention-mention pairs or mention-entity pairs), or whether learning is performed by ranking all possible candidates (mentions or entities) for a given mention.

14.2.1 Mention-pair and mention ranking models

In the **mention-pair model**, a label $y_{ij} \in \{0, 1\}$ is assigned to each pair of mentions $\langle i, j \rangle, i < j$. If i and j corefer, then $y_{ij} = 1$, and we say that i is the **antecedent** of j ; otherwise, $y_{ij} = 0$. Thus, the mention-pair model reduces coreference resolution to binary classification, enabling the application of off-the-shelf machine learning algorithms: Soon et al. (2001) use decision trees, and Bengtson and Roth (2008) use the averaged perceptron.

Under the constraint that each mention has at most one antecedent, the **antecedent structure** $\{y_{ij}\}$ induces a unique set of entities c . However, the converse is not true: a single set of entities c may be compatible with multiple antecedent structures. Since the ground truth annotations give c but not y , additional heuristics must be employed to convert the labeled data into training examples. Furthermore, we must impose the constraint that each mention have at most one antecedent. One solution is to pair the classifier with a search heuristic, based on SMASH: search backwards from j until finding an antecedent i which corefers with j with high confidence, and then stop searching. During training, for each reference j with antecedent i , we include negative examples $y_{i'j} = 0$ for all $i < i' < j$.

In **mention ranking**, the classifier learns to identify a single antecedent $a_i \in \{1, 2, \dots, i-1, i\}$ for each referring expression i , where $a_i = i$ indicates that the mention i does not refer to any previously-introduced entity (Denis and Baldridge, 2007). Specifically, the model chooses,

$$\hat{a}_i = \operatorname{argmax}_{a \in \{1, 2, \dots, i\}} \theta \cdot f(i, a, w), \quad [14.2]$$

where $f(i, a, w)$ defines a set of features on the mention pair $\langle i, a \rangle$. A special set of features can be employed for the case $a_i = i$, although later work on mention ranking has

(c) Jacob Eisenstein 2014-2017. Work in progress.

employed a two-stage model, in which an “anaphoricity” classifier determines whether the mention i refers to a previously defined entity; if so, then the ranking decision is performed over the set $1, 2, \dots, i - 1$ (Denis and Baldridge, 2008).

As with the binary coreference variables $\{y_{ij}\}$, the antecedent variables $\{a_i\}$ relate to the entity chains in a many-to-one mapping: each set of assignment variables induces a single entity clustering, but an entity clustering can correspond to many different settings of assignment variables. When mention i has multiple possible antecedents in the clustering c , a typical approach is to select the most recent compatible antecedent. However, by using a probabilistic ranking model, $p(\{a_i\} \mid \mathbf{w})$, Durrett and Klein (2013) are able to sum over the set of all antecedent structures $\mathcal{A}(c)$ that are compatible with the gold coreference clustering c ,

$$p(a_i \mid i, \mathbf{w}) = \frac{\exp(\boldsymbol{\theta} \cdot \mathbf{f}(i, a_i, \mathbf{w}))}{\sum_{a' \leq i} \exp(\boldsymbol{\theta} \cdot \mathbf{f}(i, a', \mathbf{w}))} \quad [14.3]$$

$$p(\mathbf{a} \mid \mathbf{w}) = \prod_i^M p(a_i \mid i, \mathbf{w}) \quad [14.4]$$

$$p(c \mid \mathbf{w}) = \sum_{\mathbf{a} \in \mathcal{A}(c)} p(\mathbf{a} \mid \mathbf{w}). \quad [14.5]$$

In this way, Durrett and Klein learn a model that tries to assign high scores for all valid antecedent structures.

14.2.2 Entity-based models

Many of the practical difficulties with mention-based models stem from the fact that they treat coreference resolution like a classification or ranking problem, when in fact it is a clustering problem: the goal is to group the mentions together into clusters that correspond to the underlying entities. Entity-based approaches attempt to identify these clusters directly. Such methods require defining features at the entity level, measuring whether the set of mentions are internally consistent. Cardie and Wagstaff (1999) provide an early example of entity-based coreference, incrementally merging clusters of mentions under the constraint that all pairs of mentions in the entity are compatible in number, gender, animacy, etc. They define a set of soft preferences for merging when there are multiple clusters that are compatible. More recent methods for entity-based coreference resolution have applied machine learning in the context of incremental search over the space of coreference clusterings (e.g., Clark and Manning, 2015).

The gap between entity-based and mention-pair models can be partially bridged by enforcing transitivity on the mention-pair variables: if $y_{ij} = 1$ and $y_{jk} = 1$, then $y_{ik} = 1$. This constraint can be written as a linear inequality,

$$\forall i < j < k, y_{ik} \geq y_{ij} + y_{jk} - 1. \quad [14.6]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

The transitivity constraint ensures that each mention is linked to all antecedents in its cluster. We can then formulate the inference problem as,

$$\max \sum_{i,j} \theta \cdot f(y_{ij}, w) \quad [14.7]$$

$$s.t. \forall i < j < k, y_{ik} \geq y_{ij} + y_{jk} - 1. \quad [14.8]$$

In this formulation, features are still defined over mention pairs — rather than over entire entities — but transitivity ensures that all pairs in the cluster are compatible, avoiding incoherent clusters like $\{Hillary\ Clinton \leftarrow Clinton \leftarrow Mr\ Clinton\}$. However, this coherence comes at a computational price: the constrained optimization problem is NP-hard. **Integer linear programming** (ILP), which we saw in chapter 12, is one solution (Klenner, 2007; Finkel and Manning, 2008); correlational clustering is another (McCallum and Wellner, 2004).

[todo: discuss recent methods for using deep learning to acquire entity representations (Wiseman et al., 2016; Clark and Manning, 2015)]

14.2.3 Deterministic methods

Unlike many other areas of natural language processing, it is possible to build competitive systems for coreference resolution without machine learning (Haghighi and Klein, 2009). One such architecture is shown in Figure 14.5. The basic idea is to apply a series of rule-based methods, or “sieves”, starting with high-precision rules and progressively increasing recall. Each sieve builds on the output of its predecessor, so that it is possible to consider entity-level information. For example, in the case of $\{Hillary\ Clinton \leftarrow Clinton \leftarrow she\}$, the name-matching sieve would link *Clinton* and *Hillary Clinton*, and the pronoun-matching sieve would then link *she* to the combined cluster.

The Stanford deterministic system made a strong showing at 2011 CoNLL shared task on coreference, winning nearly every track in the competition (Pradhan et al., 2011). This was particularly surprising, given the dominance of non-deterministic methods based on machine learning in virtually all other areas of natural language processing. While learning-based systems have regained the upper hand in recent years (e.g., Björkelund and Kuhn, 2014; Wiseman et al., 2016), the resurgence of deterministic, rule-based artificial intelligence in coreference resolution tells us that it may differ in important ways from other tasks, such as tagging and dependency parsing.

(c) Jacob Eisenstein 2014-2017. Work in progress.

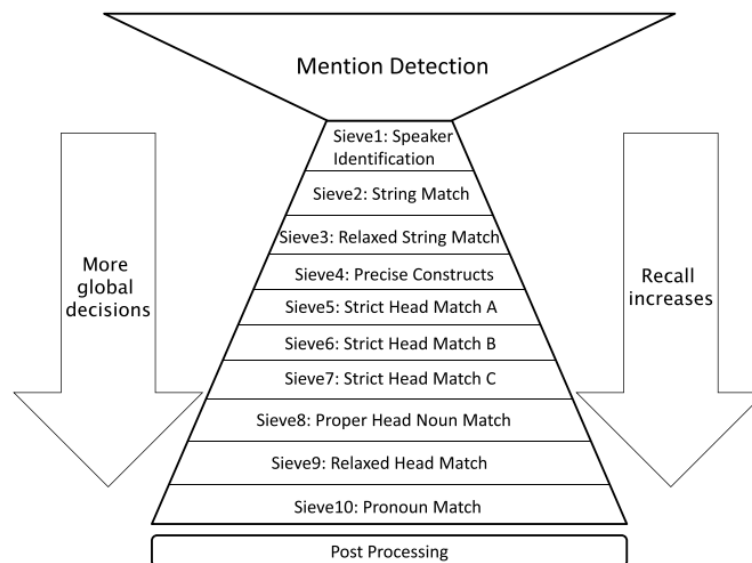


Figure 14.5: Architecture of Stanford's deterministic "multi-pass sieve" coreference system (Lee et al., 2013)

14.3 Cross-document coreference resolution

Exercises

1. The size of the largest entity typically grows linearly with the number of mentions in a document. Using this assumption, give an asymptotic estimate of the number of antecedent structures that are compatible with a coreference clustering in a document with M mentions.

Chapter 15

Discourse

15.1 Discourse relations in the Penn Discourse Treebank

- introduce discourse relations
- PDTB annotation framework in D-LTAG
- PDTB parsing

15.2 Rhetorical Structure Theory

- Higher-level discourse structure
- Shift-reduce parsing
- Applications to summarization

15.3 Centering

- Pronouns, forms of reference
- Smooth/rough transitions
- Entity grid implementation

Part IV

Applications

Chapter 16

Information extraction

Computers offer powerful capabilities for searching and reasoning about structured records and relational data. Some have even argued that the most important limitation of contemporary artificial intelligence is not inference or learning, but simply having too little knowledge (Lenat et al., 1990). Natural language processing provides an appealing solution: automatically construct a structured **knowledge base** by reading natural language text.

Many Wikipedia pages have an “infobox” that provides structured information about the an entity or event. An example is shown in Figure 16.1a: each row represents one or more properties of the entity IN THE AEROPLANE OVER THE SEA, a record album. The set of properties is determined by a predefined **schema**, which applies to all record albums in Wikipedia. As shown in Figure 16.1b, the values for many of these fields are indicated directly in the first few sentences of text on the same Wikipedia page.

The task of automatically constructing (or “populating”) an infobox based on the text is an example of **information extraction**. Much of information extraction can be described in terms of **entities**, **relations**, and **events**.

- **Entities** are uniquely specified objects in the world, such as people (Jeff Mangum), places (Athens, Georgia), organizations (Merge Records), and times (February 10, 1998). In chapter 7, we encountered the task of **named entity recognition**, which labels tokens as parts of entity spans. In information extraction, we must go further, **linking** each entity **mention** to an element in a **knowledge base**.
- **Relations** include a **predicate** and two **arguments**: for example, CAPITAL(Georgia, Atlanta).
- **Events** involve multiple typed arguments. For example, the production and release

Studio album by Neutral Milk Hotel	
Released	February 10, 1998
Recorded	July–September 1997
Studio	Pet Sounds Studio, Denver, Colorado
Genre	Indie rock • psychedelic folk • lo-fi
Length	39:55
Label	Merge • Domino
Producer	Robert Schneider

(a) A Wikipedia infobox

(16.1) In the Aeroplane Over the Sea is the second and final studio album by the American indie rock band Neutral Milk Hotel.

(16.2) It was released in the United States on February 10, 1998 on Merge Records and May 1998 on Blue Rose Records in the United Kingdom.

(16.3) Jeff Mangum moved from Athens, Georgia to Denver, Colorado to prepare the bulk of the album's material with producer Robert Schneider, this time at Schneider's newly created Pet Sounds Studio at the home of Jim McIntyre.

(b) The first few sentences of text. Strings that match fields or field names in the infobox are underlined; strings that mention other entities are wavy underlined.

Figure 16.1: From the Wikipedia page for the album “In the Aeroplane Over the Sea”, retrieved October 26, 2017.

of the album described in Figure 16.1 is described by the event,

```
<TITLE:In the Aeroplane Over the Sea,
  ARTIST:Neutral Milk Hotel,
  RELEASE-DATE:1998-Feb-10,...>
```

The set of arguments for an event type is defined by a **schema**. Events often refer to time-delimited occurrences: weddings, protests, purchases, terrorist attacks.

Information extraction is similar to predicate-argument semantic parsing tasks, such as semantic role labeling (chapter 12): we may think of predicates as corresponding to events, and the arguments as defining slots in the event representation. However, the goals of information extraction are usually more limited. Rather than accurately parsing every sentence, information extraction systems often focus on recognizing a few key relation or event types, or on the task of identifying all properties of a given entity. Information extraction is often evaluated by the correctness of the resulting knowledge base, and not by how many sentences were accurately parsed. Many relations and events will be mentioned multiple times in a corpus, but in information extraction, we are usually interested in identifying each relation and event only once — thus the goal here is sometimes described as **macro-reading**, as opposed to **micro-reading**, in which each sentence

(c) Jacob Eisenstein 2014-2017. Work in progress.

must be analyzed correctly. Macro-reading systems are not penalized for ignoring difficult sentences, as long as they can recover the same information from other, easier-to-read sources. However, macro-reading systems must resolve apparent inconsistencies (was the album released on Merge Records or Blue Rose Records?), requiring reasoning across the entire dataset.

In addition to the basic tasks of recognizing entities, relations, and events, accurate information extraction systems must handle negation, and must be able to distinguish statements of fact from hopes, fears, hunches, and hypotheticals. Finally, information extraction is often paired with the problem of **question answering**, which requires accurately parsing a query, and then selecting or generating a textual answer. Question answering systems can be built on knowledge bases that are extracted from large text corpora, or may attempt to identify answers directly from the source texts.

16.1 Entities

The starting point for information extraction is to identify mentions of entities in text. Consider the following text:

(16.4) *The United States Army captured a hill overlooking Atlanta on May 14, 1864.*

For this sentence, we have two goals:

1. **Identify** the spans *United States Army*, *Atlanta*, and *May 14, 1864* as entity mentions. (The hill is not uniquely identified, so it is not a *named* entity.) We may also want to recognize the **named entity types**: organization, location, and date. This is **named entity recognition**, and is described in chapter 7.
2. **Link** these spans to entities in a knowledge base: U.S. Army, Atlanta, and 1864–May–14. This task is known as **entity linking**.

The strings to be linked to entities are known as **mentions** — similar to the use of this term in coreference resolution.

- In some formulations of the entity linking task, only named entities are candidates for linking. This is sometimes called **named entity linking** (Ling et al., 2015c).
- In other formulations, such as **Wikification** (Milne and Witten, 2008), any string can be a mention.

The set of target entities often corresponds to Wikipedia pages, and Wikipedia is the basis for more comprehensive knowledge bases such as YAGO (Suchanek et al., 2007), DBPedia (Auer et al., 2007), and Freebase (Bollacker et al., 2008). Entity linking may also be

(c) Jacob Eisenstein 2014-2017. Work in progress.

performed in more “closed” settings, where a much smaller list of targets is provided in advance. The system must also determine if a mention does not refer to any entity in the knowledge base, sometimes called a **NIL entity** (McNamee and Dang, 2009).

Returning to (16.4), the three entity mentions may seem unambiguous. But the Wikipedia disambiguation page for the string *Atlanta* says otherwise:¹ there are more than twenty different towns and cities, five United States Navy vessels, a magazine, a television show, a band, and a singer — each prominent enough to have its own Wikipedia page. We now consider how to choose among these dozens of possibilities. In this chapter we will focus on supervised approaches. Unsupervised entity linking is closely related to the problem of **cross-document coreference resolution**, discussed in § 14.3.

16.1.1 Entity linking by learning to rank

Entity linking is often formulated as a **ranking** problem,

$$\hat{y} = \operatorname{argmax}_{y \in \mathcal{Y}(x)} \psi(y, x, c), \quad [16.1]$$

where y is a target entity, x is a description of the mention, $\mathcal{Y}(x)$ is a set of candidate entities, and c is a description of the context — such as the other text in the document, or its metadata. The function ψ is a scoring function, which could be a linear model, $\psi(y, x, c) = \theta \cdot f(y, x, c)$, or a more complex function such as a neural network. In either case, the scoring function can be learned by minimizing a margin-based **ranking loss**,

$$\ell(\hat{y}, y^{(i)}, x^{(i)}, c^{(i)}) = \left(\psi(\hat{y}, x^{(i)}, c^{(i)}) - \psi(y^{(i)}, x^{(i)}, c^{(i)}) + 1 \right)_+, \quad [16.2]$$

where $y^{(i)}$ is the ground truth and $\hat{y} \neq y^{(i)}$ is the predicted target for mention $x^{(i)}$ in context $c^{(i)}$ (Joachims, 2002; Dredze et al., 2010).

Candidate identification For computational tractability, it is helpful to restrict the set of candidates, $\mathcal{Y}(x^{(i)})$. One approach is to use a **name dictionary**, which maps from strings to the entities that they might mention. This mapping is many-to-many: a string such as *Atlanta* can refer to multiple entities, and conversely, an entity such as *Atlanta* can be referenced by multiple strings. A name dictionary can be extracted from Wikipedia, with links between each Wikipedia entity page and the anchor text of all hyperlinks that point to the page (Bunescu and Pasca, 2006; Ratnov et al., 2011). To improve recall, the name dictionary can be augmented by partial and approximate matching (Dredze et al., 2010), but as the set of candidates grows, the risk of false positives (and low precision) increases. For example, the string *Atlanta* is a partial match to *the Atlanta Fed* (a name for the Federal Reserve Bank of Atlanta), and a noisy match (edit distance of one) from *Atalanta* (a heroine in Greek mythology and an Italian soccer team).

¹[https://en.wikipedia.org/wiki/Atlanta_\(disambiguation\)](https://en.wikipedia.org/wiki/Atlanta_(disambiguation)), retrieved November 1, 2017.

Features Feature-based approaches to entity ranking rely on three main types of local information (Dredze et al., 2010):

- The similarity of the mention string to the canonical entity name, as quantified by string similarity. This feature would elevate the city *Atlanta* over the basketball team *Atlanta Hawks* for the string *Atlanta*.
- The popularity of the entity, which can be measured by Wikipedia page views or PageRank in the Wikipedia link graph. This feature would elevate *Atlanta, Georgia* over the unincorporated community of *Atlanta, Ohio*.
- The entity type, as output by the named entity recognition system. This feature would elevate the city of *Atlanta* over the magazine *Atlanta* in contexts where the mention is tagged as a location.

In addition to these local features, the document context can also help. If *Jamaica* is mentioned in a document about the Caribbean, it is likely to refer to the island nation; in the context of New York, it is likely to refer to the neighborhood in Queens; in the context of a menu, it might refer to a hibiscus tea beverage. Such hints can be formalized by computing the similarity between the Wikipedia page describing each candidate entity and the context $c^{(i)}$, which may include the bag-of-words representing the document (Dredze et al., 2010; Hoffart et al., 2011) or a smaller window of text around the mention (Ratinov et al., 2011). Contextual similarity can be modeled using the cosine similarity of bag-of-words vectors, typically weighted using **inverse document frequency** to emphasize rare words.²

Neural entity linking An alternative approach to entity ranking is to compute the score for each entity candidate using distributed vector representations of the entities, mentions, and context. For example, for the task of entity linking in Twitter, Yang et al. (2016) employ the bilinear scoring function,

$$\psi(y, x, c) = \mathbf{v}_y^\top \mathbf{W}^{(y,x)} \mathbf{x} + \mathbf{v}_y^\top \mathbf{W}^{(y,c)} \mathbf{c}, \quad [16.3]$$

with $\mathbf{v}_y \in \mathbb{R}^{K_y}$ as the vector embedding of entity y , $\mathbf{x} \in \mathbb{R}^{K_x}$ as the embedding of the mention, $\mathbf{c} \in \mathbb{R}^{K_c}$ as the embedding of the context, and the matrices $\mathbf{W}^{(y,x)}$ and $\mathbf{W}^{(y,c)}$ as parameters that score the compatibility of each entity with respect to the mention and context. Each of the vector embeddings can be learned from an end-to-end objective, or pre-trained on unlabeled data.

²The **document frequency** of word j is $\text{DF}(j) = \frac{1}{N} \sum_{i=1}^N \delta(x_j^{(i)} > 0)$, equal to the number of documents in which the word appears. The contribution of each word to the cosine similarity of two bag-of-words vectors can be weighted by the **inverse document frequency** $\frac{1}{\text{DF}(j)}$ or $\log \frac{1}{\text{DF}(j)}$, to emphasize rare words (Spärck Jones, 1972).

- Pretrained **entity embeddings** can be obtained from an existing knowledge base (Bordes et al., 2011, 2013), or by running a word embedding algorithm such as WORD2VEC on the text of wikipedia, with hyperlinks substituted for the anchor text.³
- The embedding of the mention x can be computed by averaging the embeddings of the words in the mention (Yang et al., 2016), or by one of the compositional techniques described in § 13.8.
- The embedding of the context c can be represented by a low-dimensional vector. In an auto-encoding framework, this vector is the result of noisy compression, so that it is possible to approximately reconstruct the original document (Vincent et al., 2010; Kingma and Welling, 2014). He et al. (2013) apply this idea to entity linking. The vector c can also be obtained by convolution on the embeddings of words in the document (Sun et al., 2015), or by examining metadata such as the author’s social network (Yang et al., 2016).

The remaining parameters $\mathbf{W}^{(y,x)}$ and $\mathbf{W}^{(y,c)}$ can be trained by backpropagation from the margin loss in Equation 16.2.

16.1.2 Collective entity linking

Consider the following lists:

(16.5) California, Oregon, Washington

(16.6) Baltimore, Washington, Philadelphia

(16.7) Washington, Adams, Jefferson

In each case, the term *Washington* refers to a different entity, and this reference is strongly suggested by the other entries on the list. In the last list, all three names are ambiguous in isolation — there are dozens of other *Adams* and *Jefferson* entities in Wikipedia. But a preference for coherence motivates **collectively** linking these references to the first three U.S. presidents.

A general approach to collective entity linking is to introduce a compatibility score $\psi_c(\mathbf{y})$. Collective entity linking is then performed by optimizing the global objective,

$$\hat{\mathbf{y}} = \operatorname{argmax}_{\mathbf{y} \in \mathbb{Y}(\mathbf{x})} \psi_c(\mathbf{y}) + \sum_{i=1}^N \psi_\ell(y^{(i)}, \mathbf{x}^{(i)}, \mathbf{c}^{(i)}), \quad [16.4]$$

where $\mathbb{Y}(\mathbf{x})$ is the set of all possible collective entity assignments for the mentions in \mathbf{x} , and ψ_ℓ is the local scoring function for each entity i . For simplicity, the compatibility func-

³Pre-trained entity embeddings can be downloaded from <https://code.google.com/archive/p/word2vec/>.

tion is typically decomposed into a sum of pairwise scores, $\psi_c(\mathbf{y}) = \sum_{i=1}^N \sum_{j \neq i}^N \psi_c(y^{(i)}, y^{(j)})$. The compatibility can then be computed in a number of different ways:

- Wikipedia defines high-level categories for entities (e.g., *living people*, *Presidents of the United States*, *States of the United States*), and ψ_c can reward entity pairs for the number of categories that they have in common (Cucerzan, 2007).
- Compatibility can be measured by the number of incoming hyperlinks shared by the Wikipedia pages for the two entities (Milne and Witten, 2008).
- In a neural architecture, the compatibility of two entities can be set equal to the inner product of their embeddings, $\psi_c(y^{(i)}, y^{(j)}) = \mathbf{v}_{y^{(i)}} \cdot \mathbf{v}_{y^{(j)}}$.
- A non-pairwise compatibility score can be defined using a type of latent variable model known as a **probabilistic topic model** (Blei et al., 2003; Blei, 2012). In this framework, each latent topic is a probability distribution over entities, and each document has a probability distribution over topics. Each entity helps to determine the document's distribution over topics, and in turn these topics help to resolve ambiguous entity mentions (Newman et al., 2006). Inference can be performed using the techniques described in chapter 4.

Unfortunately, collective entity linking is **NP-hard** even for pairwise compatibility functions. This means that it almost certainly cannot be solved efficiently (see Appendix B). Various approximate inference techniques have been proposed, including **integer linear programming** (Cheng and Roth, 2013), **Gibbs sampling** (Han and Sun, 2012), and graph-based algorithms (Hoffart et al., 2011; Han et al., 2011).

16.1.3 *Pairwise ranking loss functions

The loss function defined in Equation 16.2 considers only the highest-scoring prediction \hat{y} , but in fact, the true entity $y^{(i)}$ should outscore *all* other entities. A loss function based on this idea would give a gradient against the features or representations of several entities, not just the top-scoring prediction. Usunier et al. (2009) define a general ranking error function,

$$L_{\text{rank}}(k) = \sum_{j=1}^k \alpha_j, \quad \text{with } \alpha_1 \geq \alpha_2 \geq \dots \geq 0, \quad [16.5]$$

where k is equal to the number of labels ranked higher than the correct label $y^{(i)}$. This function defines a class of ranking errors: if $\alpha_j = 1$ for all j , then the ranking error is equal to the rank of the correct entity; if $\alpha_1 = 1$ and $\alpha_{j>1} = 0$, then the ranking error is one whenever the correct entity is not ranked first; if α_j decreases smoothly with j , as in $\alpha_j = \frac{1}{j}$, then the error is between these two extremes.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Algorithm 16 WARP approximate ranking loss

```

1: procedure WARP( $y^{(i)}, \mathbf{x}^{(i)}$ )
2:    $N \leftarrow 0$ 
3:   repeat
4:     Randomly sample  $y \sim \mathcal{Y}(\mathbf{x}^{(i)})$ 
5:      $N \leftarrow N + 1$ 
6:     if  $\psi(y, \mathbf{x}^{(i)}) + 1 > \psi(y^{(i)}, \mathbf{x}^{(i)})$  then ▷ check for margin violation
7:        $r \leftarrow \lfloor |\mathcal{Y}(\mathbf{x}^{(i)})|/N \rfloor$  ▷ compute approximate rank
8:       return  $L_{\text{rank}}(r) \times (\psi(y, \mathbf{x}^{(i)}) + 1 - \psi(y^{(i)}, \mathbf{x}^{(i)}))$ 
9:   until  $N \geq |\mathcal{Y}(\mathbf{x}^{(i)})| - 1$  ▷ no violation found
10:  return 0 ▷ return zero loss

```

This ranking error will be integrated into a margin objective. Remember that large margin classification requires not only the correct label, but also that the correct label outscores other labels by a substantial margin. A similar principle applies to ranking: we want a high rank for the correct entity, and we want it to be separated from other entities by a substantial margin. We therefore define the margin-augmented rank,

$$r(y^{(i)}, \mathbf{x}^{(i)}) \triangleq \sum_{y \in \mathcal{Y}(\mathbf{x}^{(i)}) \setminus y^{(i)}} \delta \left(1 + \psi(y, \mathbf{x}^{(i)}) \geq \psi(y^{(i)}, \mathbf{x}^{(i)}) \right), \quad [16.6]$$

where $\delta(\cdot)$ is a delta function, and $\mathcal{Y}(\mathbf{x}^{(i)}) \setminus y^{(i)}$ is the set of all entity candidates minus the true entity $y^{(i)}$. The margin-augmented rank is the rank of the true entity, after augmenting every other candidate with a margin of one, under the current scoring function ψ . (The context c is omitted for clarity, and can be considered part of \mathbf{x} .)

For each instance, a hinge loss is computed from the ranking error associated with this margin-augmented rank, and the violation of the margin constraint,

$$\ell(y^{(i)}, \mathbf{x}^{(i)}) = \frac{L_{\text{rank}}(r(y^{(i)}, \mathbf{x}^{(i)}))}{r(y^{(i)}, \mathbf{x}^{(i)})} \sum_{y \in \mathcal{Y}(\mathbf{x}) \setminus y^{(i)}} \left(\psi(y, \mathbf{x}^{(i)}) - \psi(y^{(i)}, \mathbf{x}^{(i)}) + 1 \right)_+, \quad [16.7]$$

The sum in Equation 16.7 includes non-zero values for every label that is ranked at least as high as the true entity, after applying the margin augmentation. Dividing by the margin-augmented rank of the true entity thus gives the average violation.

The objective in Equation 16.7 is expensive to optimize when the label space is large — as is usually the case for entity linking against large knowledge bases. This motivates a randomized approximation called **WARP** (Weston et al., 2011), shown in Algorithm 16. In this procedure, we sample random entities until one violates the pairwise margin constraint, $\psi(y, \mathbf{x}^{(i)}) + 1 \geq \psi(y^{(i)}, \mathbf{x}^{(i)})$. The number of samples N required to find

CAUSE-EFFECT	<i>those cancers were caused by radiation exposures</i>
INSTRUMENT-AGENCY	<i>phone operator</i>
PRODUCT-PRODUCER	<i>a factory manufactures suits</i>
CONTENT-CONTAINER	<i>a bottle of honey was weighed</i>
ENTITY-ORIGIN	<i>letters from foreign countries</i>
ENTITY-DESTINATION	<i>the boy went to bed</i>
COMPONENT-WHOLE	<i>my apartment has a large kitchen</i>
MEMBER-COLLECTION	<i>there are many trees in the forest</i>
COMMUNICATION-TOPIC	<i>the lecture was about semantics</i>

Table 16.1: Relations and example sentences from the SemEval-2010 dataset (Hendrickx et al., 2009)

such a violation yields an approximation of the margin-augmented rank of the true entity, $r(y^{(i)}, x^{(i)}) \approx \left\lfloor \frac{|\mathcal{Y}(x)|}{N} \right\rfloor$. If a violation is found immediately, $N = 1$, the correct entity probably ranks below many others, $r \approx |\mathcal{Y}(x)|$. If many samples are required before a violation is found, $N \rightarrow |\mathcal{Y}(x)|$, then the correct entity is probably highly ranked, $r \rightarrow 1$. A computational advantage of WARP is that it is not necessary to find the highest-scoring label, which can impose a non-trivial computational cost when $\mathcal{Y}(x^{(i)})$ is large. Note the similarity to the **negative sampling** objective in WORD2VEC (chapter 13).

16.2 Relations

Consider the following example:

(16.8) George Bush traveled to France on Thursday for a summit.

This sentence introduces a relation between the entities referenced by *George Bush* and *France*. In the Automatic Content Extraction (ACE) ontology (Linguistic Data Consortium, 2005), the type of this relation is PHYSICAL, and the subtype is LOCATED. This relation would be written,

PHYSICAL.LOCATED(*George Bush*, *France*). [16.8]

Relations take exactly two arguments, and the order of the arguments matters: the converse relation would imply that *France* is inside of *George Bush*, which is completely different.

In the ACE datasets, relations are annotated between entity mentions, as in the example above. Relations can also hold between nominals, as in the following example from the SemEval-2010 shared task (Hendrickx et al., 2009):

(c) Jacob Eisenstein 2014-2017. Work in progress.

(16.9) The cup contained tea from dried ginseng.

This sentence describes a relation of type ENTITY-ORIGIN between *tea* and *ginseng*. Nominal relation extraction is closely related to **semantic role labeling** (chapter 12). The key difference is that relation extraction is restricted to a relatively small number of relation types; for example, Table 16.1 shows the ten relation types from SemEval-2010.

16.2.1 Pattern-based relation extraction

Early work on relation extraction focused on hand-crafted patterns (Hearst, 1992). For example, the appositive *Starbuck, a native of Nantucket* signals the relation ENTITY-ORIGIN between *Starbuck* and *Nantucket*. This pattern can be written as,

PERSON , a native of LOCATION \Rightarrow ENTITY-ORIGIN(PERSON, LOCATION). [16.9]

This pattern will be “triggered” whenever the literal string *, a native of* occurs between an entity of type PERSON and an entity of type LOCATION. Such patterns can be generalized beyond literal matches using techniques such as lemmatization, which would enable the words (*buy, buys, buying*) to trigger the same patterns (see § 3.3.1). A more aggressive strategy would be to group all words in a WordNet synset (§ 3.2), grouping words like *buy* and *purchase*.

Relation extraction patterns can be implemented in finite-state automata (§ 8.1). If the named entity recognizer is also a finite-state machine, then the systems can be combined by finite-state transduction (Hobbs et al., 1997). This makes it possible to propagate uncertainty through the finite-state cascade. Suppose the entity recognizer hypothesizes that *Starbuck* refers to either a PERSON or a LOCATION; in the composed transducer, the relation extractor would be free to select the PERSON annotation when it appears in the context of an appropriate pattern.

16.2.2 Relation extraction as a classification task

Relation extraction can be formulated as a classification problem,

$$\hat{r}_{(i,j),(m,n)} = \operatorname{argmax}_{r \in \mathcal{R}} \psi(r, (i, j), (m, n), \mathbf{w}), \quad [16.10]$$

where $r \in \mathcal{R}$ is a relation type (possibly NIL), $\mathbf{w}_{i:j}$ is the span of the first argument, and $\mathbf{w}_{m:n}$ is the span of the second argument. The argument $\mathbf{w}_{m:n}$ may appear before or after $\mathbf{w}_{i:j}$ in the text, or they may overlap; we stipulate only that $\mathbf{w}_{i:j}$ is the first argument of the relation. We now consider three alternatives for computing the scoring function ψ .

(c) Jacob Eisenstein 2014-2017. Work in progress.

Feature-based classification

To build a feature-based classifier, we define the scoring function as,

$$\psi(r, (i, j), (m, n), \mathbf{w}) = \boldsymbol{\theta} \cdot \mathbf{f}(r, (i, j), (m, n), \mathbf{w}), \quad [16.11]$$

with $\boldsymbol{\theta}$ representing a vector of weights, and $\mathbf{f}(\cdot)$ a vector of features. The pattern-based methods described in § 16.2.1 suggest several features:

- Local features of $w_{i:j}$ and $w_{m:n}$, including: the strings themselves; whether they are recognized as entities, and if so, which type; whether the strings are present in a **gazetteer** of entity names; each string's syntactic **head** (§ 8.2.2).
- Features of the span between the two arguments, $w_{j:m}$ or $w_{n:i}$ (depending on which argument appears first): the length of the span; the specific words that appear in the span, either as a literal sequence or a bag-of-words; the wordnet synsets (§ 3.2) that appear in the span between the arguments.
- Features of the syntactic relationship between the two arguments, typically the **dependency path** between the arguments (§ 12.2.1). Example dependency paths are shown in Table 16.2.

Kernels

Suppose that the first line of Table 16.2 is a labeled example, and the remaining lines are instances to be classified. A feature-based approach would have to decompose the dependency paths into features that capture individual edges, with or without their labels, and then learn weights for each of these features: for example, the second line contains identical dependencies, but different arguments; the third line contains a different inflection of the word *travel*; the fourth and fifth lines each contain an additional edge on the dependency path; and the sixth example uses an entirely different path. Rather than attempting to create local features that capture all of the ways in which these dependencies paths are similar and different, we can instead define a similarity function κ , which computes a score for any pair of instances, $\kappa : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}_+$. The score for any pair of instances (i, j) is $\kappa(\mathbf{x}^{(i)}, \mathbf{x}^{(j)}) \geq 0$, with $\kappa(i, j)$ being large when instances $\mathbf{x}^{(i)}$ and $\mathbf{x}^{(j)}$ are similar. If the function κ obeys a few key properties it is called a **kernel function**.⁴

Given a valid kernel function, we can build a non-linear classifier without explicitly defining a feature vector. For a binary classification problem $y \in \{-1, 1\}$, we have the

⁴The **Gram matrix** \mathbf{K} arises from computing the kernel function between all pairs in a set of instances. For a valid kernel, the Gram matrix must be symmetric ($\mathbf{K} = \mathbf{K}^\top$) and positive semi-definite ($\forall \mathbf{a}, \mathbf{a}^\top \mathbf{K} \mathbf{a} \geq 0$). For more on kernel-based classification, see chapter 14 of Murphy (2012).

1. <i>George Bush traveled to France</i>	<i>George Bush</i> $\xleftarrow{\text{NSUBJ}}$ <i>traveled</i> $\xrightarrow{\text{OBL}}$ <i>France</i>
2. <i>Ahab traveled to Nantucket</i>	<i>Ahab</i> $\xleftarrow{\text{NSUBJ}}$ <i>traveled</i> $\xrightarrow{\text{OBL}}$ <i>Nantucket</i>
3. <i>George Bush will travel to France</i>	<i>George Bush</i> $\xleftarrow{\text{NSUBJ}}$ <i>travel</i> $\xrightarrow{\text{OBL}}$ <i>France</i>
4. <i>George Bush wants to travel to France</i>	<i>George Bush</i> $\xleftarrow{\text{NSUBJ}}$ <i>wants</i> $\xrightarrow{\text{XCOMP}}$ <i>travel</i> $\xrightarrow{\text{OBL}}$ <i>France</i>
5. <i>Ahab traveled to a city in France</i>	<i>Ahab</i> $\xleftarrow{\text{NSUBJ}}$ <i>traveled</i> $\xrightarrow{\text{OBL}}$ <i>city</i> $\xrightarrow{\text{NMOD}}$ <i>France</i>
6. <i>We await Ahab's visit to France</i>	<i>Ahab</i> $\xleftarrow{\text{NMOD:POSS}}$ <i>visit</i> $\xrightarrow{\text{NMOD}}$ <i>France</i>

Table 16.2: Candidates instances for the PHYSICAL.LOCATED relation, and their dependency paths

decision function,

$$\hat{y} = \text{Sign}(b + \sum_{i=1}^N y^{(i)} \alpha^{(i)} \kappa(\mathbf{x}^{(i)}, \mathbf{x})) \quad [16.12]$$

where b and $\{\alpha^{(i)}\}_{i=1}^N$ are parameters that must be learned from the training set, under the constraint $\forall_i, \alpha^{(i)} \geq 0$. Intuitively, each α_i specifies the importance of the instance $\mathbf{x}^{(i)}$ towards the classification rule. Kernel-based classification can be viewed as a weighted form of the **nearest-neighbor** classifier (Hastie et al., 2009), in which test instances are assigned the most common label among their near neighbors in the training set. This results in a non-linear classification boundary. The parameters are typically learned from a margin-based objective (see § 1.3), leading to the **kernel support vector machine**. To generalize to multi-class classification, we can train separate binary classifiers for each label (sometimes called **one-versus-all**), or train binary classifiers for each pair of possible labels (**one-versus-one**).

Dependency kernels are particularly effective for relation extraction, due to their ability to capture syntactic properties of the path between the two candidate arguments. One class of dependency tree kernels is defined recursively, with the score for a pair of trees equal to the similarity of the root nodes and the sum of similarities of matched pairs of child subtrees (Zelenko et al., 2003; Culotta and Sorensen, 2004). Alternatively, Bunescu and Mooney (2005) define a kernel function over sequences of unlabeled dependency edges, in which the score is computed as a product of scores for each pair of words in the sequence: identical words receive a high score, words that share a synset or part-of-speech receive a small non-zero score (e.g., *travel* / *visit*), and unrelated words receive a score of zero.

(c) Jacob Eisenstein 2014-2017. Work in progress.

Neural relation extraction

Convolutional neural networks were an early neural architecture for relation extraction (Zeng et al., 2014; dos Santos et al., 2015). For the sentence (w_1, w_2, \dots, w_M) , obtain a matrix of word embeddings \mathbf{X} , where $\mathbf{x}_m \in \mathbb{R}^K$ is the embedding of w_m . Now, suppose the candidate arguments appear at positions a_1 and a_2 ; then for each word in the sentence, its position with respect to each argument is $m - a_1$ and $m - a_2$. (Following Zeng et al. (2014), we consider a restricted version of the relation extraction task in which the arguments are single tokens.) We can augment the word embeddings by estimating embeddings for these positional offsets, $\mathbf{x}_{m-a_1}^{(p)}$ and $\mathbf{x}_{m-a_2}^{(p)}$. The complete base representation of the sentence is,

$$\mathbf{X}(a_1, a_2) = \begin{pmatrix} \mathbf{x}_1 & \mathbf{x}_2 & \cdots & \mathbf{x}_M \\ \mathbf{x}_{1-a_1}^{(p)} & \mathbf{x}_{2-a_1}^{(p)} & \cdots & \mathbf{x}_{M-a_1}^{(p)} \\ \mathbf{x}_{1-a_2}^{(p)} & \mathbf{x}_{2-a_2}^{(p)} & \cdots & \mathbf{x}_{M-a_2}^{(p)} \end{pmatrix}, \quad [16.13]$$

where each column is a vertical concatenation of a word embedding, represented by the column vector \mathbf{x}_m , and two positional embeddings, specifying the position with respect to a_1 and a_2 . The matrix $\mathbf{X}(a_1, a_2)$ is then taken as input to a convolutional layer, and max-pooling is applied to obtain a vector. The final scoring function is then,

$$\psi(r, i, j, \mathbf{X}) = \theta_r \cdot \text{MaxPool}(\text{ConvNet}(\mathbf{X}(i, j); \phi)), \quad [16.14]$$

where ϕ defines the parameters of the convolutional operator, and the θ_r defines a set of weights for relation r . The model can be trained using a margin objective,

$$\hat{r} = \underset{r}{\operatorname{argmax}} \psi(r, i, j, \mathbf{X}) \quad [16.15]$$

$$\ell = (1 + \psi(\hat{r}, i, j, \mathbf{X}) - \psi(r, i, j, \mathbf{X}))_+. \quad [16.16]$$

Recurrent neural networks have also been applied to relation extraction, using a network such as an bidirectional LSTM to encode the words or dependency path between the two arguments. Xu et al. (2015) segment each dependency path into left and right subpaths: the path,

$$(16.10) \quad \text{George Bush} \xleftarrow{\text{NSUBJ}} \text{wants} \xrightarrow{\text{XCOMP}} \text{travel} \xrightarrow{\text{OBL}} \text{France}$$

is segmented into the subpaths,

$$(16.11) \quad \text{George Bush} \xleftarrow{\text{NSUBJ}} \text{wants}$$

$$(16.12) \quad \text{wants} \xrightarrow{\text{XCOMP}} \text{travel} \xrightarrow{\text{OBL}} \text{France}.$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

They then run recurrent networks from the arguments to the root word (in this case, *wants*), obtaining the final representation by max pooling across all the recurrent states along each path. This process can be applied across separate “channels”, in which the inputs consist of embeddings for the words, parts-of-speech, dependency relations, and WordNet hypernyms. To define the model formally, let $s(m)$ define the successor of word m in either the left or right subpath (in a dependency path, each word can have a successor in at most one subpath). Let $x_m^{(c)}$ indicate the embedding of word (or relation) m in channel c , and let $\overleftarrow{h}_m^{(c)}$ indicate the associated recurrent state in the left subtree. Then the complete model is specified as follows,

$$h_{s(m)}^{(c)} = \text{RNN}(x_{s(m)}^{(c)}, h_m^{(c)}) \quad [16.17]$$

$$z^{(c)} = \text{MaxPool} \left(\overleftarrow{h}_i^{(c)}, \overleftarrow{h}_{s(i)}^{(c)}, \dots, \overleftarrow{h}_{\text{root}}^{(c)}, \overrightarrow{h}_j^{(c)}, \overrightarrow{h}_{s(j)}^{(c)}, \dots, \overrightarrow{h}_{\text{root}}^{(c)} \right) \quad [16.18]$$

$$\psi(r, i, j) = \theta \cdot \left[z^{(\text{word})}; z^{(\text{POS})}; z^{(\text{dependency})}; z^{(\text{hypernym})} \right]. \quad [16.19]$$

Note that z is computed by applying max-pooling to the *matrix* of horizontally concatenated vectors h , while ψ is computed from the *vector* of vertically concatenated vectors z . Xu et al. (2015) pass the score ψ through a **softmax** layer to obtain a probability $p(r \mid i, j, w)$, and train the model by regularized **cross-entropy**. Miwa and Bansal (2016) show that a related model can solve the more challenging “end-to-end” relation extraction task, in which the model must simultaneously detect entities and then extract their relations.

16.2.3 Knowledge base population

In many applications, what matters is not what fraction of sentences are analyzed correctly, but how much accurate knowledge can be extracted. **Knowledge base population (KBP)** refers to the task of filling in Wikipedia-style infoboxes, as shown in Figure 16.1a. Knowledge base population can be decomposed into two subtasks: **entity linking** (described in § 16.1), and **slot filling** (Ji and Grishman, 2011). Slot filling has two key differences from the formulation of relation extraction presented above:

- The relations hold between entities, rather than spans of text
- Performance is evaluated at the *type level* (on entity pairs), rather than on the *token level* (on individual sentences).

From a practical standpoint, there are three other important differences between slot filling and per-sentence relation extraction.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- KBP tasks are often formulated from the perspective of identifying attributes of a few “query” entities. As a result, these systems often start with an **information retrieval** phase, in which relevant passages of text are obtained by search.
- For many entity pairs, there will be multiple passages of text that provide evidence. Aggregating this evidence to predict a single relation type (or set of relations) is a challenge for slot filling systems.
- Labeled data is usually available in the form of pairs of related entities, rather than annotated passages of text. Training from such type-level annotations is relatively complex: two entities may be linked by several relations, or they may appear together in a passage of text that nonetheless does not describe their relation to each other.

Information retrieval is beyond the scope of this text (see Manning et al., 2008). The remainder of this section describes approaches to information fusion and learning from type-level annotations.

Information fusion

In knowledge base population, there will often be multiple pieces of evidence for (and sometimes against) a single relation. For example, a search for the entity `Maynard Jackson, Jr.` may return several passages that reference the entity `Atlanta`:

- (16.13) Elected mayor of **Atlanta** in 1973, **Maynard Jackson** was the first African American to serve as mayor of a major southern city.
- (16.14) **Atlanta**’s airport will be renamed to honor **Maynard Jackson**, the city’s first Black mayor .
- (16.15) Born in Dallas, Texas in 1938, **Maynard Holbrook Jackson, Jr.** moved to **Atlanta** when he was 8.
- (16.16) **Maynard Jackson** has gone from one of the worst high schools in **Atlanta** to one of the best.

The first and second examples provide evidence for the relation `MAYOR` holding between the entities `Atlanta` and `Maynard Jackson, Jr.`. The third example provides evidence for a different relation between these same entities, `LIVED-IN`. The fourth example poses an entity linking problem, referring to `Maynard Jackson high school`. Knowledge base population requires aggregating this sort of textual evidence, and predicting the relations that are most likely to hold.

One approach is to run a single-document relation extraction system (using the techniques described in § 16.2.2), and then aggregate the results (Li et al., 2011). Relations

that are detected with high confidence in multiple documents are more likely to be valid, motivating the heuristic,

$$\psi(r, e_1, e_2) = \sum_{i=1}^N (p(r(e_1, e_2) \mid \mathbf{w}^{(i)}))^{\alpha}, \quad [16.20]$$

where $p(r(e_1, e_2) \mid \mathbf{w}^{(i)})$ is the probability of relation r between entities e_1 and e_2 conditioned on the text $\mathbf{w}^{(i)}$, and $\alpha \gg 1$ is a tunable hyperparameter. Using this heuristic, it is possible to rank all candidate relations, and trace out a **precision-recall curve** as more relations are extracted [todo: example?]. Alternatively, features can be aggregated across multiple passages of text, feeding a single type-level relation extraction system (Wolfe et al., 2017).

Precision can be improved by introducing constraints across multiple relations. For example, if we are certain of the relation $\text{PARENT}(e_1, e_2)$, then it cannot also be the case that $\text{PARENT}(e_2, e_1)$. Integer linear programming makes it possible to incorporate such constraints into a global optimization (Li et al., 2011). Other pairs of relations have positive correlations, such as $\text{MAYOR}(e_1, e_2)$ and $\text{LIVED-IN}(e_1, e_2)$. Compatibility across relation types can be incorporated into probabilistic graphical models (e.g., Riedel et al., 2010).

Distant supervision

Relation extraction is “annotation hungry,” because each relation requires its own labeled data. Rather than relying on annotations of individual documents, it would be preferable to use existing knowledge resources — such as the many facts that are already captured in knowledge bases like DBpedia. However such annotations raise the inverse of the information fusion problem considered above: the existence of the relation $\text{MAYOR}(\text{Maynard Jackson Jr.}, \text{Atlanta})$ provides only **distant supervision** for the example texts in which this entity pair is mentioned.

One approach is to treat the entity pair as the instance, rather than the text itself (Mintz et al., 2009). Features are then aggregated across all sentences in which both entities are mentioned, and labels correspond to the relation (if any) between the entities in a knowledge base, such as FreeBase. Negative instances are constructed from entity pairs that are not related in the knowledge base. In some cases, two entities are related, but the knowledge base is missing the relation; however, because the number of possible entity pairs is huge, these missing relations are presumed to be relatively rare. This approach is shown in Figure 16.2.

In **multiple instance learning**, sets of instances receive a single label, which applies only to an unknown subset (Dietterich et al., 1997; Maron and Lozano-Pérez, 1998). This formalizes the framework of distant supervision: the relation $\text{REL}(a, b)$ can act as a label for the entire set of sentences mentioning entities a and b , even when only a subset of

(c) Jacob Eisenstein 2014-2017. Work in progress.

- **Label** : MAYOR(Atlanta, Maynard Jackson)
 - Elected mayor of **Atlanta** in 1973, **Maynard Jackson** ...
 - **Atlanta**'s airport will be renamed to honor **Maynard Jackson**, the city's first Black mayor
 - Born in Dallas, Texas in 1938, **Maynard Holbrook Jackson, Jr.** moved to **Atlanta** when he was 8.
- **Label** : MAYOR(New York, Fiorello La Guardia)
 - **Fiorello La Guardia** was Mayor of **New York** for three terms ...
 - **Fiorello La Guardia**, then serving on the **New York** City Board of Aldermen...
- **Label** : BORN-IN(Dallas, Maynard Jackson)
 - Born in **Dallas**, Texas in 1938, **Maynard Holbrook Jackson, Jr.** moved to Atlanta when he was 8.
 - **Maynard Jackson** was raised in **Dallas** ...
- **Label** : NIL(New York, Maynard Jackson)
 - **Jackson** married Valerie Richardson, whom he had met in **New York**...
 - **Jackson** was a member of the Georgia and **New York** bars ...

Figure 16.2: Four training instances for relation classification using **distant supervision** Mintz et al. (2009). The first two instances are positive for the MAYOR relation, and the third instance is positive for the BORN-IN relation. The fourth instance is a negative example, constructed from a pair of entities (New York, Maynard Jackson) that do not appear in any Freebase relation. Each instance's features are computed by aggregating across all sentences in which the two entities are mentioned.

these sentences actually describes the relation. One approach to multi-instance learning is to introduce a binary **latent variable** for each sentence, indicating whether the sentence expresses the labeled relation (Riedel et al., 2010). A variety of inference techniques have been employed for this probabilistic model of relation extraction: Surdeanu et al. (2012) use expectation maximization, Riedel et al. (2010) use sampling, and Hoffmann et al. (2011) use a custom graph-based algorithm. Expectation maximization and sampling are surveyed in chapter 4, and are covered in more detail by Murphy (2012); graph-based methods are surveyed by Mihalcea and Radev (2011).

16.2.4 Open information extraction

In classical relation extraction, the set of relations is defined in advance, using a **schema**. The relation for any pair of entities can then be predicted using multi-class classification. In **open information extraction**, a relation can be any triple of text. The example sentence (16.13) instantiates several "relations" of this sort:

(c) Jacob Eisenstein 2014-2017. Work in progress.

Task	Relation ontology	Supervision
PropBank semantic role labeling	VerbNet	sentence
FrameNet semantic role labeling	FrameNet	sentence
Relation extraction	ACE, TAC, SemEval, etc	sentence
Slot filling	ACE, TAC, SemEval, etc	relation
Open Information Extraction	open	seed relations or patterns

Table 16.3: Various relation extraction tasks and their properties. VerbNet and FrameNet are described in chapter 12. ACE (Linguistic Data Consortium, 2005), TAC (McNamee and Dang, 2009), and SemEval (Hendrickx et al., 2009) refer to shared tasks, each of which involves an ontology of relation types.

- *⟨mayor of, Maynard Jackson, Atlanta⟩*,
- *⟨elected, Maynard Jackson, mayor of Atlanta⟩*,
- *⟨elected in, Maynard Jackson, 1973⟩*,

and so on. Extracting such tuples can be viewed as a lightweight version of **semantic role labeling** (chapter 12), with only two argument types: first slot and second slot. The task is generally evaluated on the relation level, rather than on the level of sentences: precision is measured by the number of extracted relations that are accurate, and recall is measured by the number of true relations that were successfully extracted. OpenIE systems are trained from distant supervision or bootstrapping, rather than from labeled sentences.

An early example is the *TextRunner* system (Banko et al., 2007), which identifies relations with a set of handcrafted syntactic rules. The examples that are acquired from the handcrafted rules are then used to train a speedier classification model that uses part-of-speech patterns as features. Finally, the relations that are extracted by the classifier are aggregated, removing redundant relations and computing the number of times that each relation is mentioned in the corpus. *TextRunner* was the first in a series of systems that performed increasingly accurate open relation extraction by incorporating more precise linguistic features (Etzioni et al., 2011), distant supervision from Wikipedia infoboxes (Wu and Weld, 2010), and better learning algorithms (Zhu et al., 2009).

16.3 Events

Relations link pairs of entities, but many real-world situations involve more than two entities. Consider again the example sentence (16.13), which describes the **event** of an election, with four properties: the office (*mayor*), the district (*Atlanta*), the date (*1973*), and the person elected (*Maynard Jackson, Jr.*). In **event detection**, a schema is provided

(c) Jacob Eisenstein 2014-2017. Work in progress.

for each event type (e.g., an election, a terrorist attack, or a chemical reaction), indicating all the possible properties of the event. The system is then required to fill in as many of these properties as possible (Doddington et al., 2004).

Event detection systems generally involve a retrieval component (finding relevant documents and passages of text) and an extraction component (determining the properties of the event based on the retrieved texts). Early approaches focused on finite-state patterns for identify event properties (Hobbs et al., 1997); such patterns can be automatically induced by searching for patterns that are especially likely to appear in documents that match the event query (Riloff, 1996). Contemporary approaches employ techniques that are similar to FrameNet semantic role labeling (§ 12.2), such as structured prediction over local and global features (Li et al., 2013) and bidirectional recurrent neural networks (Feng et al., 2016). These methods detect whether an event is described in a sentence, and if so, what are its properties.

Event coreference Because multiple sentences may describe unique properties of a single event, **event coreference** is required to link event mentions across a single passage of text, or between passages (Humphreys et al., 1997). Bejan and Harabagiu (2014) define event coreference as the task of identifying event mentions that share the same event participants (i.e., the slot-filling entities) and the same event properties (e.g., the time and location), within or across documents. Event coreference resolution can be performed using supervised learning techniques in a similar way to entity coreference, as described in chapter 14: move left-to-right through the document, and use a classifier to decide whether to link each event reference to an existing cluster of coreferent events, or to create a new cluster (Ahn, 2006). Each clustering decision is based on the compatibility of features describing the participants and properties of the event. Due to the difficulty of annotating large amounts of data for entity coreference, unsupervised approaches are especially desirable (Chen and Ji, 2009; Bejan and Harabagiu, 2014). [todo: figure with example]

Relations between events Just as entities are related to other entities, events may be related to other events: for example, the event of winning an election both *precedes* and *causes* the event of serving as mayor; moving to Atlanta *precedes* and *enables* the event of becoming mayor of Atlanta; moving from Dallas to Atlanta *prevents* the event of later becoming mayor of Dallas. As these examples show, events may be related both temporally and causally. The **TimeML** annotation scheme specifies a set of six temporal relations between events (Pustejovsky et al., 2005), derived in part from **interval algebra** (Allen, 1984). The TimeBank corpus provides TimeML annotations for 186 documents (Pustejovsky et al., 2003). Methods for detecting these temporal relations have combined supervised machine learning with temporal constraints, such as transitivity (Mani et al., 2006; Chambers and Jurafsky, 2008).

(c) Jacob Eisenstein 2014-2017. Work in progress.

	Positive (+)	Negative (-)	Underspecified (u)
Certain (CT)	Fact: CT+	Counterfact: CT-	Certain, but unknown: CTU
Probable (PR)	Probable: PR+	Not probable: PR-	(NA)
Possible (PS)	Possible: PS+	Not possible: PS-	(NA)
Underspecified (U)	(NA)	(NA)	Unknown or uncommitted: UU

Table 16.4: Table of factuality values from the FACTBANK corpus (Saurí and Pustejovsky, 2009). The entry (NA) indicates that this combination is not annotated.

More recent annotation schemes and datasets have attempted to combine temporal and causal relations (Mirza et al., 2014; Dunietz et al., 2017): for example, the CaTeRS dataset includes annotations of 320 five-sentence short stories (Mostafazadeh et al., 2016). Abstracting still further, **processes** are networks of causal relations between multiple events. A small dataset of biological processes is annotated in the `ProcessBank` dataset (Berant et al., 2014), with the goal of supporting automatic question answering on scientific textbooks.

16.4 Hedges, denials, and hypotheticals

The methods described thus far apply to **propositions** about the way things are in the real world. But natural language can also describe events and relations that are likely or unlikely, possible or impossible, desired or feared. The following examples hint at the scope of the problem (Prabhakaran et al., 2010):

- (16.17) GM will lay off workers.
- (16.18) A spokesman for GM said GM will lay off workers.
- (16.19) GM may lay off workers.
- (16.20) The politician claimed that GM will lay off workers.
- (16.21) Some wish GM would lay off workers.
- (16.22) Will GM lay off workers?
- (16.23) Many wonder whether GM will lay off workers.

Accurate information extraction requires handling these **extra-propositional** aspects of meaning, which are sometimes summarized under the terms **modality** and **negation**.⁵

⁵The classification of negation as extra-propositional is controversial: Packard et al. (2014) argue that negation is a “core part of compositionally constructed logical-form representations.” Negation is an element of the semantic parsing tasks discussed in chapter 11 and chapter 12 — for example, negation markers are treated as adjuncts in PropBank semantic role labeling. However, many of the relation extraction methods

Modality refers to expressions of the speaker's attitude towards her own statements, including "degree of certainty, reliability, subjectivity, sources of information, and perspective" (Morante and Sporleder, 2012). Various systematizations of modality have been proposed (e.g., Palmer, 2001), including categories such as future, interrogative, imperative, conditional, and subjective. Information extraction is particularly concerned with negation and certainty. For example, Saurí and Pustejovsky (2009) link negation with a modal calculus of certainty, likelihood, and possibility, creating the two-dimensional analysis shown in Table 16.4. This schema is the basis for the FACTBANK corpus, with annotations of the **factuality** of all sentences in 208 documents of news text.

A related concept is **hedging**, in which speakers limit their commitment to a proposition (Lakoff, 1973):

(16.24) These results **suggest** that expression of c-jun, jun B and jun D genes **might** be involved in terminal granulocyte differentiation. . . (Morante and Daelemans, 2009)

(16.25) A whale is **technically** a mammal (Lakoff, 1973)

In the first example, the hedges *suggest* and *might* communicate uncertainty; in the second example, there is no uncertainty, but the hedge *technically* indicates that the evidence for the proposition will not fully meet the reader's expectations. Hedging has been studied extensively in scientific texts (Medlock and Briscoe, 2007; Morante and Daelemans, 2009), where the goal of large-scale extraction of scientific facts is obstructed by hedges and speculation. Still another related aspect of modality is **evidentiality**, in which speakers mark the source of their information. In many languages, it is obligatory to mark evidentiality through affixes or particles (Aikhenvald, 2004); while evidentiality is not grammaticalized in English, authors are expected to express this information in contexts such as journalism (Kovach and Rosenstiel, 2014) and Wikipedia.⁶

Methods for handling negation and modality generally include two phases:

1. detecting negated or uncertain events;
2. identifying the scope and focus of the negation or modal operator.

A considerable body of work on negation has employed rule-based techniques such as regular expressions (Chapman et al., 2001) to detect negated events. Such techniques match lexical cues (e.g., *Norwood was **not** elected Mayor*), while avoiding "double negatives" (e.g., *surely all this is **not without** meaning*). More recent approaches employ classi-

mentioned in this chapter do not handle negation directly. A further consideration is that negation interacts closely with aspects of modality that are generally not considered in propositional semantics, such as certainty and subjectivity.

⁶<https://en.wikipedia.org/wiki/Wikipedia:Verifiability>

fiers over lexical and syntactic features (Uzuner et al., 2009) and sequence labeling (Prabhakaran et al., 2010).

The tasks of scope and focus resolution is more fine grained, as shown in the following example from Morante and Sporleder (2012):

- (16.26) [After his habit he said] **nothing**, and after mine I asked no questions.
 After his habit he said nothing, and [after mine I asked] **no** [questions].

In this sentence, there are two negation cues (*nothing* and *no*). Each negates an event, indicated by the underlined verbs *said* and *asked* (this is the focus of negation), and each occurs within a scope: *after his habit he said* and *after mine I asked* _____ *questions*. These tasks are typically formalized as sequence labeling problems, with each word token labeled as beginning, inside, or outside of a cue, focus, or scope span (see § 7.3). Conventional sequence labeling approaches can then be applied, using surface features as well as syntax (Velldal et al., 2012) and semantic analysis (Packard et al., 2014). Labeled datasets include the BIOSCOPE corpus of biomedical texts (Vincze et al., 2008) and a shared task dataset of detective stories by Arthur Conan Doyle (Morante and Blanco, 2012).

16.5 Question answering and machine reading

The victory of the Watson question-answering system against three top human players on the game show *Jeopardy!* was a landmark moment in the history of natural language processing (Ferrucci et al., 2010). Game show questions are usually answered by **factoids**: entity names and short phrases.⁷ The task of factoid question answering is therefore closely related to information extraction, with the additional problem of accurately parsing the question.

16.5.1 Formal semantics

Semantic parsing is an effective method for question-answering in restricted domains such as questions about geography and airline reservations (Zettlemoyer and Collins, 2005), and has also been applied in “open-domain” settings such as question answering on Freebase (Berant et al., 2013) and biomedical research abstracts (Poon and Domingos, 2009). One approach is to convert the question into a lambda calculus expression that returns a boolean value: for example, the question *who is the mayor of the capital of Georgia?* would be converted to,

$$\lambda x. \exists y \text{ CAPITAL}(\text{Georgia}, y) \wedge \text{MAYOR}(y, x). \quad [16.21]$$

⁷The broader landscape of question answering includes “why” questions (*Why did Ahab continue to pursue the white whale?*), “how questions” (*How did Queequeg die?*), and requests for summaries (*What was Ishmael’s attitude towards organized religion?*). For more, see Hirschman and Gaizauskas (2001).

This lambda expression can then be used to query an existing knowledge base, returning all entities that satisfy it. The knowledge base itself could be constructed using techniques described in § 16.2.3.

16.5.2 Machine reading

Recent work has focused on answering questions about specific textual passages, similar to the reading comprehension examinations for young students (Hirschman et al., 1999). This task has come to be known as **machine reading**.

Datasets

The machine reading problem can be formulated in a number of different ways. The most important distinction is what form the answer should take.

- **Multiple-choice question answering**, as in the MCTEST dataset of stories (Richardson et al., 2013) and the New York Regents Science Exams (Clark, 2015). In MCTEST, the answer is deducible from the text alone, while in the science exams, the system must make inferences using an existing model of the underlying scientific phenomena. Here is an example from MCTEST:

(16.27) James the turtle was always getting into trouble. Sometimes he'd reach into the freezer and empty out all the food ...

Q: What is the name of the trouble making turtle?

- (a) Fries
- (b) Pudding
- (c) James
- (d) Jane

- **Cloze-style “fill in the blank” questions**, as in the CNN/DAILY MAIL comprehension task (Hermann et al., 2015), the CHILDREN'S BOOK TEST (Hill et al., 2016), and the WHO-DID-WHAT dataset (Onishi et al., 2016). In these tasks, the system must guess which word or entity completes a sentence, based on reading a passage of text. Here is an example from WHO-DID-WHAT:

(16.28) Q: Tottenham manager Juande Ramos has hinted he will allow _____ to leave if the Bulgaria striker makes it clear he is unhappy. (Onishi et al., 2016)

The query sentence may be selected either from the story itself, or from an external summary. In either case, datasets can be created automatically by processing large quantities existing documents. An additional constraint is that that missing element

(c) Jacob Eisenstein 2014-2017. Work in progress.

from the cloze must appear in the main passage of text: for example, in WHO-DID-WHAT, the candidates include all entities mentioned in the main passage. In the CNN/DAILY MAIL dataset, each entity name is replaced by a unique identifier, e.g., `entity37`. This ensures that correct answers can only be obtained by accurately reading the text, and not from external knowledge about the entities.

- **Extractive** question answering, in which the answer is drawn from the original text. In WIKIQA, answers are sentences (Yang et al., 2015); in the Stanford Question Answering Dataset (SQUAD), answers are words or short phrases (Rajpurkar et al., 2016):

(16.29) In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity.

Q: What causes precipitation to fall? A: gravity

In both WIKIQA and SQUAD, the original texts are Wikipedia articles, and the questions are generated by crowdworkers.

Methods

A baseline method is to search the text for sentences or short passages that overlap with both the query and the candidate answer (Richardson et al., 2013). In example (16.27), this baseline would select the correct answer (c), since *James* appears in a sentence that includes the query terms *trouble* and *turtle*.

This baseline can be implemented as a neural architecture, using an **attention mechanism** that scores the similarity of the query to each part of the source text (Chen et al., 2016). The first step is to encode the passage $w^{(p)}$ and the query $w^{(q)}$, using two bidirectional LSTMs (§ 6.5.4).

$$h^{(q)} = \text{BiLSTM}(w^{(q)}; \Theta^{(q)}) \quad [16.22]$$

$$h^{(p)} = \text{BiLSTM}(w^{(p)}; \Theta^{(p)}). \quad [16.23]$$

The query is represented by vertically concatenating the final states of the left-to-right and right-to-left passes:

$$u = [\overrightarrow{h^{(q)}}_{M_q}; \overleftarrow{h^{(q)}}_0]. \quad [16.24]$$

The attention vector is computed as a softmax over a vector of bilinear products, and

(c) Jacob Eisenstein 2014-2017. Work in progress.

the expected representation is computed by summing over attention values,

$$\tilde{\alpha}_m = (\mathbf{u}^{(q)})^\top \mathbf{W}_a \mathbf{h}_m^{(p)} \quad [16.25]$$

$$\boldsymbol{\alpha} = \text{SoftMax}(\tilde{\boldsymbol{\alpha}}) \quad [16.26]$$

$$\mathbf{o} = \sum_{m=1}^M \alpha_m \mathbf{h}_m^{(p)}. \quad [16.27]$$

Each candidate answer c is represented by a vector \mathbf{x}_c . Assuming the candidate answers are spans from the original text, these vectors can be set equal to the corresponding element in $\mathbf{h}^{(p)}$. The score for each candidate answer a is computed by the inner product,

$$\hat{c} = \underset{c}{\operatorname{argmax}} \mathbf{o} \cdot \mathbf{x}_c. \quad [16.28]$$

This architecture can be trained end-to-end from a loss based on the log-likelihood of the correct answer. A number of related architectures have been proposed (e.g., Hermann et al., 2015; Kadlec et al., 2016; Dhingra et al., 2017; Cui et al., 2017), and the relationships between these methods are surveyed by Wang et al. (2017).

Additional reading

The field of information extraction is surveyed in course notes by Grishman (2012), and more recently in a short survey paper (Grishman, 2015). Shen et al. (2015) survey the task of entity linking, and Ji and Grishman (2011) survey work on knowledge base population. This chapter's discussion of non-propositional meaning was strongly influenced by Morante and Sporleder (2012), who introduced a special issue of the journal *Computational Linguistics* dedicated to recent work on modality and negation.

Exercises

1. Consider the following heuristic for entity linking:

- Among all entities that have the same type as the mention (e.g., LOC, PER), choose the one whose name has the lowest edit distance from the mention.
- If more than one entity has the right type and the lowest edit distance from the mention, choose the most popular one.
- If no candidate entity has the right type, choose NIL.

Now suppose you have the following feature function:

$$\mathbf{f}(y, \mathbf{x}) = [\text{edit-dist}(\text{name}(y), \mathbf{x}), \text{same-type}(y, \mathbf{x}), \text{popularity}(y), \delta(y = \text{NIL})]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

Design a set of ranking weights θ that match the heuristic. You may assume that edit distance and popularity are always in the range $[0, 100]$, and that the NIL entity has values of zero for all features except $\delta(y = \text{NIL})$.

2. Now consider another heuristic:

- Among all candidate entities that have edit distance zero from the mention and the right type, choose the most popular one.
- If no entity has edit distance zero from the mention, choose the one with the right type that is most popular, regardless of edit distance.
- If no entity has the right type, choose NIL.

Using the same features and assumptions from the previous problem, prove that there is no set of weights that could implement this heuristic. Then show that the heuristic can be implemented by adding a single feature. Your new feature should consider only the edit distance.

3. * Consider the following formulation for collective entity linking, which rewards sets of entities that are all of the same type, where “types” can be elements of any set:

$$\psi_c(\mathbf{y}) = \begin{cases} \alpha & \text{all entities in } \mathbf{y} \text{ have the same type} \\ \beta & \text{more than half of the entities in } \mathbf{y} \text{ have the same type} \\ 0 & \text{otherwise.} \end{cases} \quad [16.29]$$

Show how to implement this model of collective entity linking in a **integer linear program**. You may want to review § 12.2.2.

To get started, here is an integer linear program for entity linking, without including the collective term ψ_c :

$$\begin{aligned} \max_{z_{i,y} \in \{0,1\}} \quad & \sum_{i=1}^N \sum_{y \in \mathcal{Y}(x^{(i)})} s_{i,y} z_{i,y} \\ \text{s.t.} \quad & \sum_{y \in \mathcal{Y}(x^{(i)})} z_{i,y} \leq 1 \quad \forall i \in \{1, 2, \dots, N\} \end{aligned}$$

where $z_{i,y} = 1$ if entity y is linked to mention i , and $s_{i,y}$ is a parameter that scores the quality of this individual ranking decision, e.g., $s_{i,y} = \theta \cdot \mathbf{f}(y, \mathbf{x}^{(i)}, \mathbf{c}^{(i)})$.

To incorporate the collective linking score, you may assume parameters r ,

$$r_{y,\tau} = \begin{cases} 1, & \text{entity } y \text{ has type } \tau \\ 0, & \text{otherwise.} \end{cases} \quad [16.30]$$

Hint: You will need to define several auxiliary variables to optimize over.

(c) Jacob Eisenstein 2014-2017. Work in progress.

4. Run `nltk.corpus.download('reuters')` to download the Reuters corpus in NLTK, and run `from nltk.corpus import reuters` to import it. The command `reuters.words()` returns an iterator over the tokens in the corpus.
- Apply the pattern `-----, such as -----` to this corpus, obtaining candidates for the IS-A relation, e.g. `IS-A(Romania, Country)`. What are three pairs that this method identifies correctly? What are three different pairs that it gets wrong?
 - Design a pattern for the PRESIDENT relation, e.g. `PRESIDENT(Philippines, Corazon Aquino)`. In this case, you may want to augment your pattern matcher with the ability to match multiple token wildcards, perhaps using case information to detect proper names. Again, list three correct
 - Preprocess the Reuters data by running a named entity recognizer, replacing tokens with named entity spans when applicable. Apply your PRESIDENT matcher to this new data. Does the accuracy improve? Compare 20 randomly-selected pairs from this pattern and the one you designed in the previous part.
5. Represent the dependency path $x^{(i)}$ as a sequence of words and dependency arcs of length M_i , ignoring the endpoints of the path. In example 1 of Table 16.2, the dependency path is,

$$x^{(1)} = (\overset{\leftarrow}{\text{NSUBJ}}, \text{traveled}, \overset{\rightarrow}{\text{OBL}}) \quad [16.31]$$

If $x_m^{(i)}$ is a word, then let $\text{pos}(x_m^{(i)})$ be its part-of-speech, using the tagset defined in chapter 7.

We can define the following kernel function over pairs of dependency paths (Bunescu and Mooney, 2005):

$$\kappa(x^{(i)}, x^{(j)}) = \begin{cases} 0, & M_i \neq M_j \\ \prod_{m=1}^{M_i} c(x_m^{(i)}, x_m^{(j)}), & M_i = M_j \end{cases}$$

$$c(x_m^{(i)}, x_m^{(j)}) = \begin{cases} 2, & x_m^{(i)} = x_m^{(j)} \\ 1, & x_m^{(i)} \text{ and } x_m^{(j)} \text{ are words and } \text{pos}(x_m^{(i)}) = \text{pos}(x_m^{(j)}) \\ 0, & \text{otherwise.} \end{cases}$$

Using this kernel function, compute the kernel similarities of example 1 from Table 16.2 with the other five examples.

6. Continuing from the previous problem, suppose that the instances have the following labels:

$$y_2 = 1, y_3 = -1, y_4 = -1, y_5 = 1, y_6 = 1 \quad [16.32]$$

Identify the conditions for α and b under which $\hat{y}_1 = 1$. Remember the constraint that $\alpha_i \geq 0$ for all i .

(c) Jacob Eisenstein 2014-2017. Work in progress.

Chapter 17

Text generation

Chapter 18

Machine translation

Appendix A

Probability

Probability theory provides a way to reason about random events. The sorts of random events that are typically used to explain probability theory include coin flips, card draws, and the weather. It may seem odd to think about the choice of a word as akin to the flip of a coin, particularly if you are the type of person to choose words carefully. But random or not, language has proven to be extremely difficult to model deterministically. Probability offers a powerful tool for modeling and manipulating linguistic data, which we will use repeatedly throughout this course.¹

Probability can be thought of in terms of **random outcomes**: for example, a single coin flip has two possible outcomes, heads or tails. The set of possible outcomes is the **sample space**, and a subset of the **sample space** is an **event**. For a sequence of two coin flips, there are four possible outcomes, $\{HH, HT, TH, TT\}$, representing the ordered sequences heads-head, heads-tails, tails-heads, and tails-tails. The event of getting exactly one head includes two outcomes: $\{HT, TH\}$.

Formally, a probability is a function from events to the interval between zero and one: $\Pr : \mathcal{F} \rightarrow [0, 1]$, where \mathcal{F} is the set of possible events. An event that is certain has probability one; an event that is impossible has probability zero. For example, the probability of getting less than three heads on two coin flips is one. Each outcome is also an event (a set with exactly one element), and for two flips of a fair coin, the probability of each outcome is,

$$\Pr(\{HH\}) = \Pr(\{HT\}) = \Pr(\{TH\}) = \Pr(\{TT\}) = \frac{1}{4}. \quad [\text{A.1}]$$

¹A good introduction to probability theory is offered by Manning and Schütze (1999), which helped to motivate this section. For more detail, Sharon Goldwater provides another useful reference, <http://homepages.inf.ed.ac.uk/sgwater/teaching/general/probability.pdf>.

A.1 Probabilities of event combinations

Because events are **sets** of outcomes, we can use set theoretic operations such as complement, intersection, and unions to reason about the probabilities of various event combinations.

For any event A , there is a **complement** $\neg A$, such that:

- The union $A \cup \neg A$ covers the entire sample space, and $\Pr(A \cup \neg A) = 1$;
- The intersection $A \cap \neg A = \emptyset$ is the empty set, and $\Pr(A \cap \neg A) = 0$.

In the coin flip example, the event of obtaining a single head on two flips corresponds to the set of outcomes $\{HT, TH\}$; the complement event includes the other two outcomes, $\{TT, HH\}$.

A.1.1 Probabilities of disjoint events

In general, when two events have an empty intersection, $A \cap B = \emptyset$, they are said to be **disjoint**. The probability of the union of two disjoint events is equal to the sum of their probabilities,

$$A \cap B = \emptyset \quad \Rightarrow \quad \Pr(A \cup B) = \Pr(A) + \Pr(B). \quad [\text{A.2}]$$

This is the **third axiom of probability**, and can be generalized to any countable sequence of disjoint events.

In the coin flip example, we can use this axiom to derive the probability of the event of getting a single head on two flips. This event is the set of outcomes $\{HT, TH\}$, which is the union of two simpler events, $\{HT, TH\} = \{HT\} \cup \{TH\}$. The events $\{HT\}$ and $\{TH\}$ are disjoint. Therefore,

$$\Pr(\{HT, TH\}) = \Pr(\{HT\} \cup \{TH\}) = \Pr(\{HT\}) + \Pr(\{TH\}) \quad [\text{A.3}]$$

$$= \frac{1}{4} + \frac{1}{4} = \frac{1}{2}. \quad [\text{A.4}]$$

For events that are not disjoint, it is still possible to compute the probability of their union:

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B). \quad [\text{A.5}]$$

This can be derived from the third axiom of probability. First, consider an event that includes all outcomes in B that are not in A , which we can write as $B - (A \cap B)$. By

(c) Jacob Eisenstein 2014-2017. Work in progress.

construction, this event is disjoint from A .² We can therefore apply the additive rule,

$$\Pr(A \cup B) = \Pr(A) + \Pr(B - (A \cap B)) \quad [\text{A.6}]$$

$$\Pr(B) = \Pr(B - (A \cap B)) + \Pr(A \cap B) \quad [\text{A.7}]$$

$$\Pr(B - (A \cap B)) = \Pr(B) - \Pr(A \cap B) \quad [\text{A.8}]$$

$$\Pr(A \cup B) = \Pr(A) + \Pr(B) - \Pr(A \cap B). \quad [\text{A.9}]$$

A.1.2 Law of total probability

A set of events $\mathcal{B} = \{B_1, B_2, \dots, B_N\}$ is a **partition** of the sample space iff each pair of events is disjoint ($B_i \cap B_j = \emptyset$), and the union of the events is the entire sample space. The law of total probability states that we can **marginalize** over these events as follows,

$$\Pr(A) = \sum_{B_n \in \mathcal{B}} \Pr(A \cap B_n). \quad [\text{A.10}]$$

Note for any event B , the union $B \cup \neg B$ forms a partition of the sample space. Therefore, an important special case of the law of total probability is,

$$\Pr(A) = \Pr(A \cap B) + \Pr(A \cap \neg B). \quad [\text{A.11}]$$

A.2 Conditional probability and Bayes' rule

A **conditional probability** is an expression like $\Pr(A \mid B)$, which is the probability of the event A , assuming that event B happens too. For example, we may be interested in 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,

$$\Pr(A \mid B) = \frac{\Pr(A \cap B)}{\Pr(B)} \quad [\text{A.12}]$$

The **chain rule** states that $\Pr(A \cap B) = \Pr(A \mid B) \times \Pr(B)$, which is just a simple rearrangement of terms from Equation A.12. We can apply the chain rule repeatedly:

$$\begin{aligned} \Pr(A \cap B \cap C) &= \Pr(A \mid B \cap C) \times \Pr(B \cap C) \\ &= \Pr(A \mid B \cap C) \times \Pr(B \mid C) \times \Pr(C) \end{aligned}$$

²[[todo: add figure](#)]

Bayes' rule (sometimes called Bayes' law or Bayes' theorem) gives us a way to convert between $\Pr(A | B)$ and $\Pr(B | A)$. It follows from the chain rule:

$$\Pr(A | B) = \frac{\Pr(A \cap B)}{\Pr(B)} = \frac{\Pr(B | A) \times \Pr(A)}{\Pr(B)} \quad [\text{A.13}]$$

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

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

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 occur on average once in 100,000 sentences. Here is an example:

(A.1) *Which class did you attend __ without registering for __?*

Lana Linguist has developed a complicated pattern matcher that attempts to identify sentences with parasitic gaps. It's pretty good, but it's not perfect:

- If a sentence has a parasitic gap, the pattern matcher will find it with probability 0.95. (Skipping ahead, 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**. The **precision** is defined as one minus the false positive rate.)

Suppose that Lana's pattern matcher says that a sentence contains a parasitic gap. What is the probability that this is true?

Let G be the event of a sentence having a parasitic gap, and T be the event of the test being positive. We are interested in the probability of a sentence having a parasitic gap given that the test is positive. This is the conditional probability $\Pr(G | T)$, and we can compute it from Bayes' rule:

$$\Pr(G | T) = \frac{\Pr(T | G) \times \Pr(G)}{\Pr(T)}. \quad [\text{A.14}]$$

(c) Jacob Eisenstein 2014-2017. Work in progress.

We already know both terms in the numerator: $\Pr(T \mid G)$ is the recall, which is 0.95; $\Pr(G)$ is the prior, which is 10^{-5} .

We are not given the denominator, but we can compute it by using some of the tools that we have developed in this section. We first apply the law of total probability, using the partition $\{G, \neg G\}$:

$$\Pr(T) = \Pr(T \cap G) + \Pr(T \cap \neg G). \quad [\text{A.15}]$$

This says that the probability of the test being positive is the sum of the probability of a **true positive** ($T \cap G$) and the probability of a **false positive** ($T \cap \neg G$). Next, we can compute the probability of each of these events using the chain rule:

$$\Pr(T \cap G) = \Pr(T \mid G) \times \Pr(G) = 0.95 \times 10^{-5} \quad [\text{A.16}]$$

$$\Pr(T \cap \neg G) = \Pr(T \mid \neg G) \times \Pr(\neg G) = 0.005 \times (1 - 10^{-5}) \approx 0.005 \quad [\text{A.17}]$$

$$\Pr(T) = \Pr(T \cap G) + \Pr(T \cap \neg G) \quad [\text{A.18}]$$

$$= 0.95 \times 10^{-5} + 0.005 \approx 0.005. \quad [\text{A.19}]$$

We now return to Bayes' rule to compute the desired posterior probability,

$$\Pr(G \mid T) = \frac{\Pr(T \mid G) \Pr(G)}{\Pr(T)} \quad [\text{A.20}]$$

$$= \frac{0.95 \times 10^{-5}}{0.95 \times 10^{-5} + 0.005 \times (1 - 10^{-5})} \quad [\text{A.21}]$$

$$\approx 0.002. \quad [\text{A.22}]$$

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

A.3 Independence

Two events are independent if the probability of their intersection is equal to the product of their probabilities: $\Pr(A \cap B) = \Pr(A) \times \Pr(B)$. For example, for two flips of a fair coin, the probability of getting heads on the first flip is independent of the probability of getting heads on the second flip. We can prove this by using the additive axiom defined

(c) Jacob Eisenstein 2014-2017. Work in progress.

above:

$$\Pr(\{HT, HH\}) = \Pr(HT) + \Pr(HH) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad [\text{A.23}]$$

$$\Pr(\{HH, TH\}) = \Pr(HH) + \Pr(TH) = \frac{1}{4} + \frac{1}{4} = \frac{1}{2} \quad [\text{A.24}]$$

$$\Pr(\{HT, HH\}) \times \Pr(\{HH, TH\}) = \frac{1}{2} \times \frac{1}{2} = \frac{1}{4} \quad [\text{A.25}]$$

$$\Pr(\{HT, HH\} \cap \{HH, TH\}) = \Pr(HH) = \frac{1}{4} \quad [\text{A.26}]$$

$$= \Pr(\{HT, HH\}) \times \Pr(\{HH, TH\}). \quad [\text{A.27}]$$

Independence will play a key role in the discussion of probabilistic classification later in this chapter.

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

A.4 Random variables

Random variables are functions of events. Formally, we will treat random variables as functions from events to the space \mathbb{R}^n , where \mathbb{R} is the set of real numbers. This general notion subsumes a number of different types of random variables:

- **Indicator random variables** are functions from events to the set $\{0, 1\}$. In the coin flip example, we can define Y as an indicator random variable, for whether the coin has come up heads on at least one flip. This would include the outcomes $\{HH, HT, TH\}$. The event probability $\Pr(Y = 1)$ is the sum of the probabilities of these outcomes, $\Pr(Y = 1) = \frac{1}{4} + \frac{1}{4} + \frac{1}{4} = \frac{3}{4}$.
- A **discrete random variable** is a function from events to a countable subset of \mathbb{R} . Consider the coin flip example: the number of heads, X , can be viewed as a discrete random variable, $X \in 0, 1, 2$. The event probability $\Pr(X = 1)$ can again be computed as the sum of the probabilities of the events in which there is one head, $\{HT, TH\}$, giving $\Pr(X = 1) = \frac{1}{2}$.

Each possible value of a random variable is associated with a subset of the sample space. In the coin flip example, $X = 0$ is associated with the event $\{TT\}$, $X = 1$ is associated with the event $\{HT, TH\}$, and $X = 2$ is associated with the event $\{HH\}$. Assuming a fair coin, the probabilities of these events are, respectively, $1/4$, $1/2$, and $1/4$. This list of numbers represents the **probability distribution** over X , written p_X , which maps from the possible values of X to the non-negative reals. For a specific value x , we

write $p_X(x)$, which is equal to the event probability $\Pr(X = x)$.³ The function p_X is called a probability **mass** function (pmf) if X is discrete; it is called a probability **density** function (pdf) if X is continuous. In either case, we have $\int_x p_X(x)dx = 1$ and $\forall x, p_X(x) \geq 0$.

Random variables can be combined into **joint probabilities**, e.g., $p_{A,B}(a, b) = \Pr(A = a \cap B = b)$. Several ideas from event probabilities carry over to probability distributions over random variables:

- We can write a **marginal probability distribution** $p_A(a) = \sum_b p_{A,B}(a, b)$.
- We can write a **conditional probability distribution** as $p_{A|B}(a | b) = \frac{p_{A,B}(a, b)}{p_B(b)}$.
- Random variables A and B are independent iff $p_{A,B}(a, b) = p_A(a) \times p_B(b)$.

A.5 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:

$$E[g(x)] = \int_{\mathcal{X}} g(x)p(x)dx \quad [\text{A.28}]$$

For example, a fast food restaurant in Quebec has a special offer for cold days: they give a 1% discount on poutine for every degree below zero. Assuming they use a thermometer with infinite precision, the expected price would be an integral over all possible temperatures,

$$E[\text{price}(x)] = \int_{\mathcal{X}} \min(1, 1 + x) \times \text{original-price} \times p(x)dx. \quad [\text{A.29}]$$

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

³In general, capital letters (e.g., X) refer to random variables, and lower-case letters (e.g., x) refer to specific values. I will often just write $p(x)$, when the subscript is clear from context.

A.6 Modeling and estimation

Probabilistic models give us a principled way to reason about random events and random variables, and to make predictions about the future. Let's consider the coin toss example. We can model each toss as a random event, with probability θ of the event H , and probability $1 - \theta$ of the complementary event T . If we write a random variable X as the total number of heads on three coin flips, then the distribution of X depends on θ . In this case, X is distributed as a **binomial random variable**, meaning that it is drawn from a binomial distribution, with **parameters** $(\theta, N = 3)$. We write:

$$X \sim \text{Binomial}(\theta, N = 3). \quad [\text{A.30}]$$

This is a probabilistic model of X . The binomial distribution has a number of known properties that enable us to make statements about the X , such as its expected value, the likelihood that its value will fall within some interval, etc.

Now suppose that θ is unknown, but we have run an experiment, in which we executed N trials, and obtained x heads. We can **estimate** θ by the principle of **maximum likelihood**:

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p_X(x; \theta, N). \quad [\text{A.31}]$$

This says that our estimate $\hat{\theta}$ should be the value that maximizes the likelihood of the data we have observed. The semicolon indicates that θ and N are parameters of the probability function. The likelihood $p_X(x; \theta, N)$ can be computed from the binomial distribution,

$$p_X(x; \theta, N) = \frac{N!}{x!(N-x)!} \theta^x (1 - \theta)^{N-x}. \quad [\text{A.32}]$$

This likelihood is proportional to the product of the probability of individual outcomes: for example, the sequence T, H, H, T, H would have probability $\theta^2(1 - \theta)^3$. The term $\frac{N!}{x!(N-x)!}$ arises from the many possible orderings by which we could obtain x heads on N trials. This term is constant in θ , so it can be ignored.

We can maximize likelihood by taking the derivative and setting it equal to zero. In practice, we usually maximize log-likelihood, which is a monotonic function of the likeli-

(c) Jacob Eisenstein 2014-2017. Work in progress.

hood, and is easier to manipulate mathematically.

$$\ell(\theta) = x \log \theta + (N - x) \log(1 - \theta) \quad [\text{A.33}]$$

$$\frac{\partial \ell(\theta)}{\partial \theta} = \frac{x}{\theta} - \frac{N - x}{1 - \theta} \quad [\text{A.34}]$$

$$\frac{N - x}{1 - \theta} = \frac{x}{\theta} \quad [\text{A.35}]$$

$$\frac{N - x}{x} = \frac{1 - \theta}{\theta} \quad [\text{A.36}]$$

$$\frac{N}{x} - 1 = \frac{1}{\theta} - 1 \quad [\text{A.37}]$$

$$\hat{\theta} = \frac{x}{N}. \quad [\text{A.38}]$$

In this case, the maximum likelihood estimate is equal to $\frac{x}{N}$, the fraction of trials that came up heads. This intuitive solution is also known as the **relative frequency estimate**, since it is equal to the relative frequency of the outcome.

Is maximum likelihood estimation always the right choice? Suppose you conduct one trial, and get heads — would you conclude that $\theta = 1$, so this coin is guaranteed to give heads? If not, then you must have some **prior expectation** about θ . To incorporate this prior information, we can treat θ as a random variable, and use Bayes rule:

$$p(\theta | x; N) = \frac{p(x | \theta) \times p(\theta)}{p(x)} \quad [\text{A.39}]$$

$$\propto p(x | \theta) \times p(\theta) \quad [\text{A.40}]$$

$$\hat{\theta} = \underset{\theta}{\operatorname{argmax}} p(x | \theta) \times p(\theta). \quad [\text{A.41}]$$

This is the **maximum a posteriori** (MAP) estimate. Given a form for $p(\theta)$, you can derive the MAP estimate using the same approach that was used to derive the maximum likelihood estimate.

Appendix B

Computational complexity

[todo: 2-3 paragraphs about computational complexity: big-O notation, P vs NP.] The reader is encouraged to consult Arora and Barak (2009) and Sipser (2012) for more on this topic.

Bibliography

- Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I. J., Harp, A., Irving, G., Isard, M., Jia, Y., Józefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D. G., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P. A., Vanhoucke, V., Vasudevan, V., Viégas, F. B., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2016). Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *CoRR*, abs/1603.04467.
- Abend, O. and Rappoport, A. (2017). The state of the art in semantic representation. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Abney, S., Schapire, R. E., and Singer, Y. (1999). Boosting applied to tagging and PP attachment. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 132–134.
- Abney, S. P. (1987). *The English noun phrase in its sentential aspect*. PhD thesis, Massachusetts Institute of Technology.
- Abney, S. P. and Johnson, M. (1991). Memory requirements and local ambiguities of parsing strategies. *Journal of Psycholinguistic Research*, 20(3):233–250.
- Ahn, D. (2006). The stages of event extraction. In *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, pages 1–8. Association for Computational Linguistics.
- Aikhenvald, A. Y. (2004). *Evidentiality*. Oxford University Press.
- Akaike, H. (1974). A new look at the statistical model identification. *IEEE Transactions on Automatic Control*, 19(6):716–723.
- Akmajian, A., Demers, R. A., Farmer, A. K., and Harnish, R. M. (2010). *Linguistics: An introduction to language and communication*. MIT press, Cambridge, MA, sixth edition.

- Allauzen, C., Riley, M., Schalkwyk, J., Skut, W., and Mohri, M. (2007). OpenFst: A general and efficient weighted finite-state transducer library. In *International Conference on Implementation and Application of Automata*, pages 11–23. Springer.
- Allen, J. F. (1984). Towards a general theory of action and time. *Artificial intelligence*, 23(2):123–154.
- Alm, C. O., Roth, D., and Sproat, R. (2005). Emotions from text: machine learning for text-based emotion prediction. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 579–586.
- Aluísio, S., Pelizzoni, J., Marchi, A., de Oliveira, L., Manenti, R., and Marquiafável, V. (2003). An account of the challenge of tagging a reference corpus for Brazilian Portuguese. *Computational Processing of the Portuguese Language*, pages 194–194.
- Anand, P., Walker, M., Abbott, R., Fox Tree, J. E., Bowmani, R., and Minor, M. (2011). Cats rule and dogs drool!: Classifying stance in online debate. In *Proceedings of the 2nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis*, pages 1–9, Portland, Oregon. Association for Computational Linguistics.
- Anandkumar, A. and Ge, R. (2016). Efficient approaches for escaping higher order saddle points in non-convex optimization. In *Proceedings of the Conference On Learning Theory (COLT)*, pages 81–102.
- 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.
- Ando, R. K. and Zhang, T. (2005). A framework for learning predictive structures from multiple tasks and unlabeled data. *The Journal of Machine Learning Research*, 6:1817–1853.
- Andor, D., Alberti, C., Weiss, D., Severyn, A., Presta, A., Ganchev, K., Petrov, S., and Collins, M. (2016). Globally normalized transition-based neural networks. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 2442–2452.
- Aronoff, M. (1976). *Word formation in generative grammar*. MIT Press.
- Arora, S. and Barak, B. (2009). *Computational complexity: a modern approach*. Cambridge University Press.
- 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.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Arora, S., Li, Y., Liang, Y., Ma, T., and Risteski, A. (2016). Linear algebraic structure of word senses, with applications to polysemy. *arXiv preprint arXiv:1601.03764*.
- Artstein, R. and Poesio, M. (2008). Inter-coder agreement for computational linguistics. *Computational Linguistics*, 34(4):555–596.
- Artzi, Y. and Zettlemoyer, L. (2013). Weakly supervised learning of semantic parsers for mapping instructions to actions. *Transactions of the Association for Computational Linguistics*, 1:49–62.
- Auer, P. (2013). *Code-switching in conversation: Language, interaction and identity*. Routledge.
- Auer, S., Bizer, C., Kobilarov, G., Lehmann, J., Cyganiak, R., and Ives, Z. (2007). Dbpedia: A nucleus for a web of open data. *The semantic web*, pages 722–735.
- Austin, J. L. (1962). *How to do things with words*. Oxford University Press.
- Aw, A., Zhang, M., Xiao, J., and Su, J. (2006). A phrase-based statistical model for SMS text normalization. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 33–40.
- Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016). Layer normalization. *arXiv preprint arXiv:1607.06450*.
- Baldwin, T. and Kim, S. N. (2010). Multiword expressions. In *Handbook of natural language processing*, volume 2, pages 267–292. CRC Press, Boca Raton, USA.
- Balle, B., Quattoni, A., and Carreras, X. (2011). A spectral learning algorithm for finite state transducers. In *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML)*, pages 156–171.
- Banarescu, L., Bonial, C., Cai, S., Georgescu, M., Griffitt, K., Hermjakob, U., Knight, K., Koehn, P., Palmer, M., and Schneider, N. (2013). Abstract meaning representation for sembanking. In *Proceedings of the 7th Linguistic Annotation Workshop and Interoperability with Discourse*, pages 178–186, Sofia, Bulgaria. Association for Computational Linguistics.
- Banko, M., Cafarella, M. J., Soderland, S., Broadhead, M., and Etzioni, O. (2007). Open information extraction from the web. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 2670–2676.
- Barber, D. (2012). *Bayesian reasoning and machine learning*. Cambridge University Press.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Barman, U., Das, A., Wagner, J., and Foster, J. (2014). Code mixing: A challenge for language identification in the language of social media. In *Proceedings of the First Workshop on Computational Approaches to Code Switching*, pages 13–23, Doha, Qatar. Association for Computational Linguistics.
- Barnickel, T., Weston, J., Collobert, R., Mewes, H.-W., and Stümpflen, V. (2009). Large scale application of neural network based semantic role labeling for automated relation extraction from biomedical texts. *PLoS One*, 4(7):e6393.
- Baron, A. and Rayson, P. (2008). Vard2: A tool for dealing with spelling variation in historical corpora. In *Postgraduate conference in corpus linguistics*.
- Baroni, M., Bernardi, R., and Zamparelli, R. (2014). Frege in space: A program for compositional distributional semantics. *Linguistic Issues in Language Technologies*.
- Baroni, M. and Zamparelli, R. (2010). Nouns are vectors, adjectives are matrices: Representing adjective-noun constructions in semantic space. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1183–1193.
- Beesley, K. R. and Karttunen, L. (2003). *Finite-state morphology*. Center for the Study of Language and Information, Stanford, CA.
- Bejan, C. A. and Harabagiu, S. (2014). Unsupervised event coreference resolution. *Computational Linguistics*, 40(2):311–347.
- 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.
- Bengio, Y., Simard, P., and Frasconi, P. (1994). Learning long-term dependencies with gradient descent is difficult. *IEEE Transactions on Neural Networks*, 5(2):157–166.
- Bengtson, E. and Roth, D. (2008). Understanding the value of features for coreference resolution. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 294–303.
- Benjamini, Y. and Hochberg, Y. (1995). Controlling the false discovery rate: a practical and powerful approach to multiple testing. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 289–300.
- Berant, J., Chou, A., Frostig, R., and Liang, P. (2013). Semantic parsing on freebase from question-answer pairs. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1533–1544.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Berant, J., Srikumar, V., Chen, P.-C., Vander Linden, A., Harding, B., Huang, B., Clark, P., and Manning, C. D. (2014). Modeling biological processes for reading comprehension. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- 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.
- Berg-Kirkpatrick, T., Burkett, D., and Klein, D. (2012). An empirical investigation of statistical significance in NLP. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 995–1005.
- 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.
- Bergstra, J., Breuleux, O., Bastien, F., Lamblin, P., Pascanu, R., Desjardins, G., Turian, J., Warde-Farley, D., and Bengio, Y. (2010). Theano: A CPU and GPU math compiler in Python. In *Proceedings of the 9th Python in Science Conference*, pages 1–7.
- Bertsekas, D. P. (2012). Incremental gradient, subgradient, and proximal methods for convex optimization: A survey. In Sra et al. (2012).
- Bhatia, P., Guthrie, R., and Eisenstein, J. (2016). Morphological priors for probabilistic neural word embeddings. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Biber, D. (1991). *Variation across speech and writing*. Cambridge University Press.
- Bikel, D. M. (2004). Intricacies of Collins’ parsing model. *Computational Linguistics*, 30(4):479–511.
- Bishop, C. M. (2006). *Pattern recognition and machine learning*. springer.
- Björkelund, A. and Kuhn, J. (2014). Learning structured perceptrons for coreference resolution with latent antecedents and non-local features. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 47–57.
- Blackburn, P. and Bos, J. (2005). *Representation and inference for natural language: A first course in computational semantics*. CSLI.
- Blei, D. M. (2012). Probabilistic topic models. *Communications of the ACM*, 55(4):77–84.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Blei, D. M. (2014). Build, compute, critique, repeat: Data analysis with latent variable models. *Annual Review of Statistics and Its Application*, 1:203–232.
- Blei, D. M., Ng, A. Y., and Jordan, M. I. (2003). Latent dirichlet allocation. *the Journal of machine Learning research*, 3:993–1022.
- 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.
- Blum, A. and Mitchell, T. (1998). Combining labeled and unlabeled data with co-training. In *Proceedings of the Conference On Learning Theory (COLT)*, pages 92–100.
- Böhmová, A., Hajič, J., Hajičová, E., and Hladká, B. (2003). The Prague dependency tree-bank. In *Treebanks, Text, Speech and Language Technology*, pages 103–127. Springer.
- Bojanowski, P., Grave, E., Joulin, A., and Mikolov, T. (2017). Enriching word vectors with subword information. *Transactions of the Association for Computational Linguistics*, 5:135–146.
- Bollacker, K., Evans, C., Paritosh, P., Sturge, T., and Taylor, J. (2008). Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the ACM International Conference on Management of Data (SIGMOD)*, pages 1247–1250. AcM.
- Bontcheva, K., Dimitrov, M., Maynard, D., Tablan, V., and Cunningham, H. (2002). Shallow methods for named entity coreference resolution. In *Chaines de références et résolveurs d’anaphores, workshop TALN*.
- Bordes, A., Usunier, N., Garcia-Duran, A., Weston, J., and Yakhnenko, O. (2013). Translating embeddings for modeling multi-relational data. In *Neural Information Processing Systems (NIPS)*, pages 2787–2795.
- Bordes, A., Weston, J., Collobert, R., Bengio, Y., et al. (2011). Learning structured embeddings of knowledge bases. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 301–306.
- Botha, J. A. and Blunsom, P. (2014). Compositional morphology for word representations and language modelling. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Bottou, L. (1998). Online learning and stochastic approximations. *On-line learning in neural networks*, pages 9–42.
- Bottou, L. (2012). Stochastic gradient descent tricks. In *Neural networks: Tricks of the trade*, pages 421–436. Springer.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Bottou, L., Curtis, F. E., and Nocedal, J. (2016). Optimization methods for large-scale machine learning. *arXiv preprint arXiv:1606.04838*.
- Bowman, S. R., Vilnis, L., Vinyals, O., Dai, A., Jozefowicz, R., and Bengio, S. (2016). Generating sentences from a continuous space. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 10–21.
- Boyd, S. and Vandenberghe, L. (2004). *Convex Optimization*. Cambridge University Press, New York.
- Branavan, S., Chen, H., Eisenstein, J., and Barzilay, R. (2009). Learning document-level semantic properties from free-text annotations. *Journal of Artificial Intelligence Research*, 34(2):569–603.
- Brown, P. F., Desouza, P. V., Mercer, R. L., Pietra, V. J. D., and Lai, J. C. (1992). Class-based n-gram models of natural language. *Computational linguistics*, 18(4):467–479.
- Brun, C. and Roux, C. (2014). Décomposition des “hash tags” pour l’amélioration de la classification en polarité des “tweets”. *Proceedings of Traitement Automatique des Langues Naturelles*, pages 473–478.
- Bruni, E., Tran, N.-K., and Baroni, M. (2014). Multimodal distributional semantics. *Journal of Artificial Intelligence Research*, 49(2014):1–47.
- Bullinaria, J. A. and Levy, J. P. (2007). Extracting semantic representations from word co-occurrence statistics: A computational study. *Behavior research methods*, 39(3):510–526.
- Bunescu, R. C. and Mooney, R. J. (2005). A shortest path dependency kernel for relation extraction. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 724–731.
- Bunescu, R. C. and Pasca, M. (2006). Using encyclopedic knowledge for named entity disambiguation. In *Proceedings of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 9–16.
- Cai, Q. and Yates, A. (2013). Large-scale semantic parsing via schema matching and lexicon extension. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 423–433.
- Canny, J. (1987). A computational approach to edge detection. In *Readings in Computer Vision*, pages 184–203. Elsevier.
- Cappé, O. and Moulines, E. (2009). On-line expectation–maximization algorithm for latent data models. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)*, 71(3):593–613.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Cardie, C. and Wagstaff, K. (1999). Noun phrase coreference as clustering. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 82–89.
- Carletta, J. (1996). Assessing agreement on classification tasks: the kappa statistic. *Computational linguistics*, 22(2):249–254.
- Carpenter, B. (1997). *Type-logical semantics*. MIT Press, Cambridge, MA.
- Carreras, X., Collins, M., and Koo, T. (2008). Tag, dynamic programming, and the perceptron for efficient, feature-rich parsing. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 9–16.
- Carreras, X. and Màrquez, L. (2005). Introduction to the conll-2005 shared task: Semantic role labeling. In *Proceedings of the Ninth Conference on Computational Natural Language Learning*, pages 152–164. Association for Computational Linguistics.
- Carroll, L. (1917). *Through the looking glass: And what Alice found there*. Rand, McNally, Chicago.
- Chambers, N. and Jurafsky, D. (2008). Jointly combining implicit constraints improves temporal ordering. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 698–706.
- Chang, M.-W., Ratinov, L., and Roth, D. (2007). Guiding semi-supervision with constraint-driven learning. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 280–287.
- Chang, M.-W., Ratinov, L.-A., Rizzolo, N., and Roth, D. (2008). Learning and inference with constraints. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 1513–1518.
- Chapman, W. W., Bridewell, W., Hanbury, P., Cooper, G. F., and Buchanan, B. G. (2001). A simple algorithm for identifying negated findings and diseases in discharge summaries. *Journal of biomedical informatics*, 34(5):301–310.
- Charniak, E. (1997). Statistical techniques for natural language parsing. *AI magazine*, 18(4):33–43.
- Charniak, E. and Johnson, M. (2005). Coarse-to-fine n-best parsing and maxent discriminative reranking. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 173–180.
- Chelba, C. and Acero, A. (2006). Adaptation of maximum entropy capitalizer: Little data can help a lot. *Computer Speech & Language*, 20(4):382–399.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Chen, D., Bolton, J., and Manning, C. D. (2016). A thorough examination of the CNN/-Daily Mail reading comprehension task. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Chen, D. and Manning, C. D. (2014). A fast and accurate dependency parser using neural networks. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 740–750.
- Chen, M., Xu, Z., Weinberger, K., and Sha, F. (2012). Marginalized denoising autoencoders for domain adaptation. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Chen, S. F. and Goodman, J. (1999). An empirical study of smoothing techniques for language modeling. *Computer Speech & Language*, 13(4):359–393.
- Chen, T. and Guestrin, C. (2016). Xgboost: A scalable tree boosting system. In *Proceedings of Knowledge Discovery and Data Mining (KDD)*, pages 785–794.
- Chen, X., Qiu, X., Zhu, C., Liu, P., and Huang, X. (2015). Long short-term memory neural networks for chinese word segmentation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1197–1206.
- Chen, Z. and Ji, H. (2009). Graph-based event coreference resolution. In *Proceedings of the 2009 Workshop on Graph-based Methods for Natural Language Processing*, pages 54–57. Association for Computational Linguistics.
- Cheng, X. and Roth, D. (2013). Relational inference for wikification. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1787–1796.
- Chiang, D., Graehl, J., Knight, K., Pauls, A., and Ravi, S. (2010). Bayesian inference for finite-state transducers. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 447–455.
- Cho, K., Van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using rnn encoder-decoder for statistical machine translation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Chomsky, N. (1957). *Syntactic structures*. Mouton & Co., The Hague.
- Choromanska, A., Henaff, M., Mathieu, M., Arous, G. B., and LeCun, Y. (2015). The loss surfaces of multilayer networks. In *Proceedings of Artificial Intelligence and Statistics (AISTATS)*, pages 192–204.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Christensen, J., Soderland, S., Etzioni, O., et al. (2010). Semantic role labeling for open information extraction. In *Proceedings of the Workshop on Formalisms and Methodology for Learning by Reading*, pages 52–60. Association for Computational Linguistics.
- Chu, Y.-J. and Liu, T.-H. (1965). On shortest arborescence of a directed graph. *Scientia Sinica*, 14(10):1396–1400.
- Chung, C. and Pennebaker, J. W. (2007). The psychological functions of function words. In Fiedler, K., editor, *Social communication*, pages 343–359. Psychology Press, New York and Hove.
- Church, K. W. (2000). Empirical estimates of adaptation: the chance of two Noriegas is closer to $p/2$ than p^2 . In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 180–186.
- Church, K. W. and Hanks, P. (1990). Word association norms, mutual information, and lexicography. *Computational linguistics*, 16(1):22–29.
- Ciaramita, M. and Johnson, M. (2003). Supersense tagging of unknown nouns in wordnet. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 168–175.
- Clark, K. and Manning, C. D. (2015). Entity-centric coreference resolution with model stacking. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1405–1415.
- Clark, P. (2015). Elementary school science and math tests as a driver for ai: take the aristo challenge! In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 4019–4021.
- Clarke, J., Goldwasser, D., Chang, M.-W., and Roth, D. (2010). Driving semantic parsing from the world’s response. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 18–27.
- Clarke, J. and Lapata, M. (2008). Global inference for sentence compression: An integer linear programming approach. *Journal of Artificial Intelligence Research*, 31:399–429.
- Cohen, J. (1960). A coefficient of agreement for nominal scales. *Educational and psychological measurement*, 20(1):37–46.
- Cohen, S. (2016). *Bayesian analysis in natural language processing*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers, San Rafael, CA.
- Cohen, S. B., Gómez-Rodríguez, C., and Satta, G. (2012). Elimination of spurious ambiguity in transition-based dependency parsing. *CoRR*, abs/1206.6735.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Collier, N., Nobata, C., and Tsujii, J.-i. (2000). Extracting the names of genes and gene products with a hidden markov model. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 201–207.
- 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*, pages 27–38.
- Collins, M. and Koo, T. (2005). Discriminative reranking for natural language parsing. *Computational Linguistics*, 31(1):25–70.
- Collobert, R., Kavukcuoglu, K., and Farabet, C. (2011a). Torch7: A matlab-like environment for machine learning. Technical Report EPFL-CONF-192376, EPFL.
- 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.
- Collobert, R., Weston, J., Bottou, L., Karlen, M., Kavukcuoglu, K., and Kuksa, P. (2011b). Natural language processing (almost) from scratch. *Journal of Machine Learning Research*, 12:2493–2537.
- Conneau, A., Kiela, D., Schwenk, H., Barrault, L., and Bordes, A. (2017). Supervised learning of universal sentence representations from natural language inference data. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 681–691.
- Cormen, T. H., Leiserson, C. E., Rivest, R. L., and Stein, C. (2009). *Introduction to algorithms*. MIT press, third edition.
- Cotterell, R., Schütze, H., and Eisner, J. (2016). Morphological smoothing and extrapolation of word embeddings. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1651–1660.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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. and Singer, Y. (2001). Pranking with ranking. In *Neural Information Processing Systems (NIPS)*, pages 641–647.
- Creutz, M. and Lagus, K. (2002). Unsupervised discovery of morphemes. In *Proceedings of the ACL-02 workshop on Morphological and phonological learning-Volume 6*, pages 21–30. Association for Computational Linguistics.
- Cucerzan, S. (2007). Large-scale named entity disambiguation based on wikipedia data. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- 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.
- Cui, Y., Chen, Z., Wei, S., Wang, S., Liu, T., and Hu, G. (2017). Attention-over-attention neural networks for reading comprehension. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Culotta, A. and Sorensen, J. (2004). Dependency tree kernels for relation extraction. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Curry, H. B. and Feys, R. (1958). *Combinatory Logic*, volume I. North Holland, Amsterdam.
- Danescu-Niculescu-Mizil, C., Sudhof, M., Jurafsky, D., Leskovec, J., and Potts, C. (2013). A computational approach to politeness with application to social factors. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 250–259.
- Das, D., Chen, D., Martins, A. F., Schneider, N., and Smith, N. A. (2014). Frame-semantic parsing. *Computational Linguistics*, 40(1):9–56.
- Daumé III, H. (2007). Frustratingly easy domain adaptation. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., and Bengio, Y. (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. In *Neural Information Processing Systems (NIPS)*, pages 2933–2941.
- Davidson, D. (1967). The logical form of action sentences. In Rescher, N., editor, *The Logic of Decision and Action*. University of Pittsburgh Press, Pittsburgh.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- De Marneffe, M.-C., Dozat, T., Silveira, N., Haverinen, K., Ginter, F., Nivre, J., and Manning, C. D. (2014). Universal stanford dependencies: A cross-linguistic typology. In *LREC*, pages 4585–4592.
- De Marneffe, M.-C., MacCartney, B., and Manning, C. D. (2006). Generating typed dependency parses from phrase structure parses. In *LREC*, pages 449–454.
- 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.
- Dean, J. and Ghemawat, S. (2008). Mapreduce: simplified data processing on large clusters. *Communications of the ACM*, 51(1):107–113.
- 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.
- Dehdari, J. (2014). *A Neurophysiologically-Inspired Statistical Language Model*. PhD thesis, The Ohio State University.
- Dempster, A. P., Laird, N. M., and Rubin, D. B. (1977). Maximum likelihood from incomplete data via the em algorithm. *Journal of the Royal Statistical Society. Series B (Methodological)*, pages 1–38.
- Denis, P. and Baldridge, J. (2007). A ranking approach to pronoun resolution. In *IJCAI*.
- Denis, P. and Baldridge, J. (2008). Specialized models and ranking for coreference resolution. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing, EMNLP '08*, pages 660–669, Stroudsburg, PA, USA. Association for Computational Linguistics.
- Dhingra, B., Liu, H., Yang, Z., Cohen, W. W., and Salakhutdinov, R. (2017). Gated-attention readers for text comprehension. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Dietterich, T. G. (1998). Approximate statistical tests for comparing supervised classification learning algorithms. *Neural computation*, 10(7):1895–1923.
- Dietterich, T. G., Lathrop, R. H., and Lozano-Pérez, T. (1997). Solving the multiple instance problem with axis-parallel rectangles. *Artificial intelligence*, 89(1):31–71.
- Dimitrova, L., Ide, N., Petkevic, V., Erjavec, T., Kaalep, H. J., and Tufis, D. (1998). Multext-east: Parallel and comparable corpora and lexicons for six central and eastern european languages. In *Proceedings of the 17th international conference on Computational linguistics-Volume 1*, pages 315–319. Association for Computational Linguistics.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Doddington, G. R., Mitchell, A., Przybocki, M. A., Ramshaw, L. A., Strassel, S., and Weischedel, R. M. (2004). The automatic content extraction (ace) program-tasks, data, and evaluation. In *Proceedings of the Language Resources and Evaluation Conference*, pages 837–840.
- dos Santos, C., Xiang, B., and Zhou, B. (2015). Classifying relations by ranking with convolutional neural networks. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 626–634.
- Dowty, D. (1991). Thematic proto-roles and argument selection. *Language*, pages 547–619.
- Dredze, M., McNamee, P., Rao, D., Gerber, A., and Finin, T. (2010). Entity disambiguation for knowledge base population. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 277–285. Association for Computational Linguistics.
- Dredze, M., Paul, M. J., Bergsma, S., and Tran, H. (2013). Carmen: A Twitter geolocation system with applications to public health. In *AAAI workshop on expanding the boundaries of health informatics using AI (HIAI)*, pages 20–24.
- 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.
- Dunietz, J., Levin, L., and Carbonell, J. (2017). The because corpus 2.0: Annotating causality and overlapping relations. In *Proceedings of the Linguistic Annotation Workshop*.
- 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*.
- Durrett, G. and Klein, D. (2015). Neural crf parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Dyer, C., Ballesteros, M., Ling, W., Matthews, A., and Smith, N. A. (2015). Transition-based dependency parsing with stack long short-term memory. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 334–343.
- Edmonds, J. (1967). Optimum branchings. *Journal of Research of the National Bureau of Standards B*, 71(4):233–240.
- Efron, B. and Tibshirani, R. J. (1993). An introduction to the bootstrap: Monographs on statistics and applied probability, vol. 57. *New York and London: Chapman and Hall/CRC*.
- Eisner, J. (2002). Parameter estimation for probabilistic finite-state transducers. In *acl*, pages 1–8.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Eisner, J. M. (1996). Three new probabilistic models for dependency parsing: An exploration. In *COLING*, pages 340–345.
- Ekman, P. (1992). Are there basic emotions? *Psychological Review*, 99(3):550–553.
- Esuli, A. and Sebastiani, F. (2006). Sentiwordnet: A publicly available lexical resource for opinion mining. In *LREC*, volume 6, pages 417–422. Citeseer.
- Etzioni, O., Fader, A., Christensen, J., Soderland, S., and Mausam, M. (2011). Open information extraction: The second generation. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, pages 3–10.
- Faruqui, M., Dodge, J., Jauhar, S. K., Dyer, C., Hovy, E., and Smith, N. A. (2015). Retrofitting word vectors to semantic lexicons. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Faruqui, M. and Dyer, C. (2014). Improving vector space word representations using multilingual correlation. In *Proceedings of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 462–471.
- Faruqui, M., McDonald, R., and Soricut, R. (2016a). Morpho-syntactic lexicon generation using graph-based semi-supervised learning. *Transactions of the Association for Computational Linguistics*, 4:1–16.
- Faruqui, M., Tsvetkov, Y., Rastogi, P., and Dyer, C. (2016b). Problems with evaluation of word embeddings using word similarity tasks. In *Proceedings of the 1st Workshop on Evaluating Vector-Space Representations for NLP*, pages 30–35, Berlin, Germany. Association for Computational Linguistics.
- Fellbaum, C. (2010). *WordNet*. Springer.
- Feng, X., Huang, L., Tang, D., Ji, H., Qin, B., and Liu, T. (2016). A language-independent neural network for event detection. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 66–71.
- Ferrucci, D., Brown, E., Chu-Carroll, J., Fan, J., Gondek, D., Kalyanpur, A. A., Lally, A., Murdock, J. W., Nyberg, E., Prager, J., et al. (2010). Building Watson: An overview of the DeepQA project. *AI magazine*, 31(3):59–79.
- 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 Bach, E. and Harms, R., editors, *Universals in linguistic theory*. Holt, Rinehart, and Winston.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Fillmore, C. J. (1976). Frame semantics and the nature of language. *Annals of the New York Academy of Sciences*, 280(1):20–32.
- Fillmore, C. J. and Baker, C. (2009). A frames approach to semantic analysis. In *The Oxford Handbook of Linguistic Analysis*. Oxford University Press.
- 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., Grenager, T., and Manning, C. D. (2007). The infinite tree. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 272–279.
- 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.
- Finkel, J. R. and Manning, C. (2009). Hierarchical bayesian domain adaptation. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 602–610.
- Finkel, J. R. and Manning, C. D. (2008). Enforcing transitivity in coreference resolution. In *Proceedings of the 46th Annual Meeting of the Association for Computational Linguistics on Human Language Technologies: Short Papers*, pages 45–48. Association for Computational Linguistics.
- Finkelstein, L., Gabrilovich, E., Matias, Y., Rivlin, E., Solan, Z., Wolfman, G., and Ruppin, E. (2002). Placing search in context: The concept revisited. *ACM Transactions on Information Systems*, 20(1):116–131.
- Firth, J. R. (1957). *Papers in Linguistics 1934-1951*. Oxford University Press.
- Flanigan, J., Thomson, S., Carbonell, J., Dyer, C., and Smith, N. A. (2014). A discriminative graph-based parser for the abstract meaning representation. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1426–1436.
- Francis, W. N. (1964). A standard sample of present-day English for use with digital computers. Report to the U.S Office of Education on Cooperative Research Project No. E-007.
- Freund, Y. and Schapire, R. E. (1999). Large margin classification using the perceptron algorithm. *Machine learning*, 37(3):277–296.
- Fromkin, V., Rodman, R., and Hyams, N. (2013). *An introduction to language*. Cengage Learning.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Fundel, K., Küffner, R., and Zimmer, R. (2007). Relex – relation extraction using dependency parse trees. *Bioinformatics*, 23(3):365–371.
- Gabrilovich, E. and Markovitch, S. (2007). Computing semantic relatedness using wikipedia-based explicit semantic analysis. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, volume 7, pages 1606–1611.
- Gage, P. (1994). A new algorithm for data compression. *The C Users Journal*, 12(2):23–38.
- Gale, W. A., Church, K. W., and Yarowsky, D. (1992). One sense per discourse. In *Proceedings of the workshop on Speech and Natural Language*, pages 233–237. Association for Computational Linguistics.
- Ganchev, K. and Dredze, M. (2008). Small statistical models by random feature mixing. In *Proceedings of the ACL08 HLT Workshop on Mobile Language Processing*, pages 19–20.
- Ganchev, K., Graça, J., Gillenwater, J., and Taskar, B. (2010). Posterior regularization for structured latent variable models. *The Journal of Machine Learning Research*, 11:2001–2049.
- Ganin, Y., Ustinova, E., Ajakan, H., Germain, P., Larochelle, H., Laviolette, F., Marchand, M., and Lempitsky, V. (2016). Domain-adversarial training of neural networks. *Journal of Machine Learning Research*, 17(59):1–35.
- 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.
- Ge, D., Jiang, X., and Ye, Y. (2011). A note on the complexity of l_p minimization. *Mathematical programming*, 129(2):285–299.
- Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015). Escaping from saddle points — online stochastic gradient for tensor decomposition. In Grünwald, P., Hazan, E., and Kale, S., editors, *Proceedings of the Conference On Learning Theory (COLT)*.
- Ge, R. and Mooney, R. J. (2005). A statistical semantic parser that integrates syntax and semantics. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 9–16.
- Geach, P. T. (1962). *Reference and generality: An examination of some medieval and modern theories*. Cornell University Press.
- Gildea, D. and Jurafsky, D. (2002). Automatic labeling of semantic roles. *Computational linguistics*, 28(3):245–288.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Gimpel, K., Schneider, N., O'Connor, B., Das, D., Mills, D., Eisenstein, J., Heilman, M., Yogatama, D., Flanigan, J., and Smith, N. A. (2011). Part-of-speech tagging for Twitter: annotation, features, and experiments. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 42–47.
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of Artificial Intelligence and Statistics (AISTATS)*, pages 249–256.
- Glorot, X., Bordes, A., and Bengio, Y. (2011). Deep sparse rectifier networks. In *Proceedings of the 14th International Conference on Artificial Intelligence and Statistics. JMLR W&CP Volume*, volume 15, pages 315–323.
- Godfrey, J. J., Holliman, E. C., and McDaniel, J. (1992). Switchboard: Telephone speech corpus for research and development. In *Proceedings of the International Conference on Acoustics, Speech, and Signal Processing (ICASSP)*, pages 517–520. IEEE.
- Goldberg, Y. (2017a). An adversarial review of “adversarial generation of natural language”. <https://medium.com/@yoav.goldberg/an-adversarial-review-of-adversarial-generation-of-natural-language-409ac3378bd7>.
- Goldberg, Y. (2017b). *Neural Network Methods for Natural Language Processing*. Synthesis Lectures on Human Language Technologies. Morgan & Claypool Publishers.
- 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.
- Goodfellow, I., Bengio, Y., and Courville, A. (2016). *Deep learning*. MIT Press.
- Goodman, J. T. (2001). A bit of progress in language modeling. *Computer Speech & Language*, 15(4):403–434.
- Gouws, S., Metzler, D., Cai, C., and Hovy, E. (2011). Contextual bearing on linguistic variation in social media. In *LASM*.
- Graves, A. and Schmidhuber, J. (2005). Framewise phoneme classification with bidirectional lstm and other neural network architectures. *Neural Networks*, 18(5):602–610.
- Grishman, R. (2012). Information extraction: Capabilities and challenges. Notes prepared for the 2012 International Winter School in Language and Speech Technologies, Rovira i Virgili University, Tarragona, Spain.
- Grishman, R. (2015). Information extraction. *IEEE Intelligent Systems*, 30(5):8–15.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Grishman, R., Macleod, C., and Sterling, J. (1992). Evaluating parsing strategies using standardized parse files. In *Proceedings of the third conference on Applied natural language processing*, pages 156–161. Association for Computational Linguistics.
- Grishman, R. and Sundheim, B. (1996). Message understanding conference-6: A brief history. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, volume 96, pages 466–471.
- Groenendijk, J. and Stokhof, M. (1991). Dynamic predicate logic. *Linguistics and philosophy*, 14(1):39–100.
- 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.
- 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. (2009). Simple coreference resolution with rich syntactic and semantic features. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1152–1161.
- Hajič, J. and Hladká, B. (1998). Tagging inflective languages: Prediction of morphological categories for a rich, structured tagset. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 483–490.
- Hammerton, J. (2003). Named entity recognition with long short-term memory. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 172–175.
- Han, X. and Sun, L. (2012). An entity-topic model for entity linking. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 105–115.
- Han, X., Sun, L., and Zhao, J. (2011). Collective entity linking in web text: a graph-based method. In *Proceedings of ACM SIGIR conference on Research and development in information retrieval*, pages 765–774.
- 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)*.
- Haspelmath, M. and Sims, A. (2013). *Understanding morphology*. Routledge.
- Hastie, T., Tibshirani, R., and Friedman, J. (2009). *The elements of statistical learning*. Springer, New York, second edition.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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.
- Hayes, A. F. and Krippendorff, K. (2007). Answering the call for a standard reliability measure for coding data. *Communication methods and measures*, 1(1):77–89.
- He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level performance on imagenet classification. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 1026–1034.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 770–778.
- He, L., Lee, K., Lewis, M., and Zettlemoyer, L. (2017). Deep semantic role labeling: What works and what’s next. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- He, Z., Liu, S., Li, M., Zhou, M., Zhang, L., and Wang, H. (2013). Learning entity representation for entity disambiguation. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 30–34.
- Hearst, M. A. (1992). Automatic acquisition of hyponyms from large text corpora. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 539–545. Association for Computational Linguistics.
- Hendrickx, I., Kim, S. N., Kozareva, Z., Nakov, P., Ó Séaghdha, D., Padó, S., Pennacchiotti, M., Romano, L., and Szpakowicz, S. (2009). Semeval-2010 task 8: Multi-way classification of semantic relations between pairs of nominals. In *Proceedings of the Workshop on Semantic Evaluations: Recent Achievements and Future Directions*, pages 94–99. Association for Computational Linguistics.
- Hermann, K. M., Kocisky, T., Grefenstette, E., Espeholt, L., Kay, W., Suleyman, M., and Blunsom, P. (2015). Teaching machines to read and comprehend. In *Advances in Neural Information Processing Systems*, pages 1693–1701.
- Hill, F., Bordes, A., Chopra, S., and Weston, J. (2016). The goldilocks principle: Reading children’s books with explicit memory representations. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Hindle, D. and Rooth, M. (1990). Structural ambiguity and lexical relations. In *Proceedings of the Workshop on Speech and Natural Language*.
- Hirschman, L. and Gaizauskas, R. (2001). Natural language question answering: the view from here. *natural language engineering*, 7(4):275–300.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Hirschman, L., Light, M., Breck, E., and Burger, J. D. (1999). Deep read: A reading comprehension system. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 325–332.
- Hobbs, J. R. (1978). Resolving pronoun references. *Lingua*, 44(4):311–338.
- Hobbs, J. R., Appelt, D., Bear, J., Israel, D., Kameyama, M., Stickel, M., and Tyson, M. (1997). Fastus: A cascaded finite-state transducer for extracting information from natural-language text. *Finite-state language processing*, pages 383–406.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural computation*, 9(8):1735–1780.
- Hockenmaier, J. and Steedman, M. (2007). Cggbank: a corpus of ccg derivations and dependency structures extracted from the penn treebank. *Computational Linguistics*, 33(3):355–396.
- Hoffart, J., Yosef, M. A., Bordino, I., Fürstenau, H., Pinkal, M., Spaniol, M., Taneva, B., Thater, S., and Weikum, G. (2011). Robust disambiguation of named entities in text. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 782–792.
- Hoffmann, R., Zhang, C., Ling, X., Zettlemoyer, L., and Weld, D. S. (2011). Knowledge-based weak supervision for information extraction of overlapping relations. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 541–550.
- Holmstrom, L. and Koistinen, P. (1992). Using additive noise in back-propagation training. *IEEE Transactions on Neural Networks*, 3(1):24–38.
- Hovy, E. and Lavid, J. (2010). Towards a ‘science’ of corpus annotation: a new methodological challenge for corpus linguistics. *International journal of translation*, 22(1):13–36.
- Hovy, E., Marcus, M., Palmer, M., Ramshaw, L., and Weischedel, R. (2006). Ontonotes: the 90% solution. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 57–60.
- 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, F. and Yates, A. (2012). Biased representation learning for domain adaptation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1313–1323.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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.
- Huang, Z., Xu, W., and Yu, K. (2015). Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.
- Huffman, D. A. (1952). A method for the construction of minimum-redundancy codes. *Proceedings of the IRE*, 40(9):1098–1101.
- Humphreys, K., Gaizauskas, R., and Azzam, S. (1997). Event coreference for information extraction. In *Proceedings of a Workshop on Operational Factors in Practical, Robust Anaphora Resolution for Unrestricted Texts*, pages 75–81. Association for Computational Linguistics.
- Ide, N. and Wilks, Y. (2006). Making sense about sense. In *Word sense disambiguation*, pages 47–73. Springer.
- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 448–456.
- Iyyer, M., Manjunatha, V., Boyd-Graber, J., and Daumé III, H. (2015). Deep unordered composition rivals syntactic methods for text classification. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1681–1691.
- Jeong, M., Lin, C.-Y., and Lee, G. G. (2009). Semi-supervised speech act recognition in emails and forums. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1250–1259.
- Ji, H. and Grishman, R. (2011). Knowledge base population: Successful approaches and challenges. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1148–1158.
- Ji, Y. and Eisenstein, J. (2015). One vector is not enough: Entity-augmented distributional semantics for discourse relations. *Transactions of the Association for Computational Linguistics (TACL)*.
- Ji, Y., Haffari, G., and Eisenstein, J. (2016). A latent variable recurrent neural network for discourse relation language models. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Ji, Y., Tan, C., Martschat, S., Choi, Y., and Smith, N. A. (2017). Dynamic entity representations in neural language models. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1831–1840.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Jiang, L., Yu, M., Zhou, M., Liu, X., and Zhao, T. (2011). Target-dependent twitter sentiment classification. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 151–160.
- Jing, H. (2000). Sentence reduction for automatic text summarization. In *Proceedings of the sixth conference on Applied natural language processing*, pages 310–315. Association for Computational Linguistics.
- Joachims, T. (2002). Optimizing search engines using clickthrough data. In *Proceedings of Knowledge Discovery and Data Mining (KDD)*, pages 133–142.
- Jockers, M. L. (2015). Szuzhet? <http://bla.bla.com>.
- Johnson, M. (1998). Pcfg models of linguistic tree representations. *Computational Linguistics*, 24(4):613–632.
- Johnson, R. and Zhang, T. (2017). Deep pyramid convolutional neural networks for text categorization. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 562–570.
- Joshi, A. K. (1985). How much context-sensitivity is required to provide reasonable structural descriptions? – tree adjoining grammars. In *Natural Language Processing – Theoretical, Computational and Psychological Perspective*. Cambridge University Press, New York, NY.
- Joshi, A. K. and Schabes, Y. (1997). Tree-adjoining grammars. In *Handbook of formal languages*, pages 69–123. Springer.
- Joshi, A. K., Shanker, K. V., and Weir, D. (1991). The convergence of mildly context-sensitive grammar formalisms. In *Foundational Issues in Natural Language Processing*. MIT Press, Cambridge MA.
- Jozefowicz, R., Zaremba, W., and Sutskever, I. (2015). An empirical exploration of recurrent network architectures. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 2342–2350.
- Jurafsky, D. (1996). A probabilistic model of lexical and syntactic access and disambiguation. *Cognitive Science*, 20(2):137–194.
- Jurafsky, D. and Martin, J. H. (2009). *Speech and Language Processing*. Prentice Hall, 2 edition.
- Kadlec, R., Schmid, M., Bajgar, O., and Kleindienst, J. (2016). Text understanding with the attention sum reader network. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 908–918.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Kalchbrenner, N. and Blunsom, P. (2013). Recurrent convolutional neural networks for discourse compositionality. In *Proceedings of the Workshop on Continuous Vector Space Models and their Compositionality*, pages 119–126, Sofia, Bulgaria. Association for Computational Linguistics.
- Kalchbrenner, N., Grefenstette, E., and Blunsom, P. (2014). A convolutional neural network for modelling sentences. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 655–665.
- Karlsson, F. (2007). Constraints on multiple center-embedding of clauses. *Journal of Linguistics*, 43(02):365–392.
- Kate, R. J., Wong, Y. W., and Mooney, R. J. (2005). Learning to transform natural to formal languages. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*.
- Kehler, A. (2007). Rethinking the SMASH approach to pronoun interpretation. In *Interdisciplinary perspectives on reference processing*, New Directions in Cognitive Science Series, pages 95–122. Oxford University Press.
- Kilgarrieff, A. (1997). I don’t believe in word senses. *Computers and the Humanities*, 31(2):91–113.
- Kim, M.-J. (2002). Does korean have adjectives? *MIT Working Papers in Linguistics*, 43:71–89.
- Kim, S.-M. and Hovy, E. (2006). Extracting opinions, opinion holders, and topics expressed in online news media text. In *Proceedings of the Workshop on Sentiment and Subjectivity in Text*, pages 1–8, Sydney, Australia. Association for Computational Linguistics.
- Kim, Y. (2014). Convolutional neural networks for sentence classification. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1746–1751.
- Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2016). Character-aware neural language models. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*.
- Kingma, D. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Kingma, D. P. and Welling, M. (2014). Auto-encoding variational bayes. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Kipper-Schuler, K. (2005). *VerbNet: A broad-coverage, comprehensive verb lexicon*. PhD thesis, Computer and Information Science, University of Pennsylvania.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Kiros, R., Salakhutdinov, R., and Zemel, R. (2014). Multimodal neural language models. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 595–603.
- Kiros, R., Zhu, Y., Salakhutdinov, R., Zemel, R. S., Torralba, A., Urtasun, R., and Fidler, S. (2015). Skip-thought vectors. In *Neural Information Processing Systems (NIPS)*.
- Klein, D. and Manning, C. D. (2003). Accurate unlexicalized parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 423–430.
- Klementiev, A., Titov, I., and Bhattarai, B. (2012). Inducing crosslingual distributed representations of words. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 1459–1474.
- Klenner, M. (2007). Enforcing consistency on coreference sets. In *Recent Advances in Natural Language Processing (RANLP)*, pages 323–328.
- 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. (2005). Hidden-variable models for discriminative reranking. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 507–514.
- Koo, T. and Collins, M. (2010). Efficient third-order dependency parsers. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1–11.
- 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.
- Kovach, B. and Rosenstiel, T. (2014). *The elements of journalism: What newspeople should know and the public should expect*. Three Rivers Press.
- Krishnamurthy, J. (2016). Probabilistic models for learning a semantic parser lexicon. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 606–616.
- Krishnamurthy, J. and Mitchell, T. M. (2012). Weakly supervised training of semantic parsers. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 754–765.
- Kübler, S., McDonald, R., and Nivre, J. (2009). Dependency parsing. *Synthesis Lectures on Human Language Technologies*, 1(1):1–127.

(c) Jacob Eisenstein 2014–2017. Work in progress.

- Kuhlmann, M. and Nivre, J. (2010). Transition-based techniques for non-projective dependency parsing. *Northern European Journal of Language Technology (NEJLT)*, 2(1):1–19.
- Kummerfeld, J. K., Berg-Kirkpatrick, T., and Klein, D. (2015). An empirical analysis of optimization for max-margin NLP. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Kwiatkowski, T., Goldwater, S., Zettlemoyer, L., and Steedman, M. (2012). A probabilistic model of syntactic and semantic acquisition from child-directed utterances and their meanings. In *Proceedings of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 234–244.
- Lafferty, J., McCallum, A., and Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *icml*.
- Lakoff, G. (1973). Hedges: A study in meaning criteria and the logic of fuzzy concepts. *Journal of philosophical logic*, 2(4):458–508.
- Lample, G., Ballesteros, M., Subramanian, S., Kawakami, K., and Dyer, C. (2016). Neural architectures for named entity recognition. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 260–270.
- Lappin, S. and Leass, H. J. (1994). An algorithm for pronominal anaphora resolution. *Computational linguistics*, 20(4):535–561.
- Lari, K. and Young, S. J. (1990). The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer speech & language*, 4(1):35–56.
- Law, E. and Ahn, L. v. (2011). Human computation. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 5(3):1–121.
- LeCun, Y. and Bengio, Y. (1995). Convolutional networks for images, speech, and time series. *The handbook of brain theory and neural networks*, 3361.
- LeCun, Y., Bottou, L., Orr, G. B., and Müller, K.-R. (1998). Efficient backprop. In *Neural networks: Tricks of the trade*, pages 9–50. Springer.
- Lee, C. M. and Narayanan, S. S. (2005). Toward detecting emotions in spoken dialogs. *IEEE transactions on speech and audio processing*, 13(2):293–303.
- Lee, H., Chang, A., Peirsman, Y., Chambers, N., Surdeanu, M., and Jurafsky, D. (2013). Deterministic coreference resolution based on entity-centric, precision-ranked rules. *Computational Linguistics*, 39(4):885–916.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Lee, H., Peirsman, Y., Chang, A., Chambers, N., Surdeanu, M., and Jurafsky, D. (2011). Stanford's multi-pass sieve coreference resolution system at the conll-2011 shared task. In *Proceedings of the Conference on Natural Language Learning (CoNLL)*, pages 28–34. Association for Computational Linguistics.
- Lenat, D. B., Guha, R. V., Pittman, K., Pratt, D., and Shepherd, M. (1990). Cyc: toward programs with common sense. *Communications of the ACM*, 33(8):30–49.
- 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). Dependency-based word embeddings. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 302–308.
- Levy, O., Goldberg, Y., and Dagan, I. (2015). Improving distributional similarity with lessons learned from word embeddings. *Transactions of the Association for Computational Linguistics*, 3:211–225.
- Lewis, M. and Steedman, M. (2013). Combined distributional and logical semantics. *Transactions of the Association for Computational Linguistics*, 1:179–192.
- Li, J. and Jurafsky, D. (2015). Do multi-sense embeddings improve natural language understanding? In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1722–1732.
- Li, J., Luong, M.-T., and Jurafsky, D. (2015a). A hierarchical neural autoencoder for paragraphs and documents. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Li, J., Luong, T., Jurafsky, D., and Hovy, E. (2015b). When are tree structures necessary for deep learning of representations? In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 2304–2314.
- Li, Q., Anzaroot, S., Lin, W.-P., Li, X., and Ji, H. (2011). Joint inference for cross-document information extraction. In *Proceedings of the International Conference on Information and Knowledge Management (CIKM)*, pages 2225–2228.
- Li, Q., Ji, H., and Huang, L. (2013). Joint event extraction via structured prediction with global features. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 73–82.
- Liang, P., Bouchard-Côté, A., Klein, D., and Taskar, B. (2006). An end-to-end discriminative approach to machine translation. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 761–768.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Liang, P., Jordan, M. I., and Klein, D. (2013). Learning dependency-based compositional semantics. *Computational Linguistics*, 39(2):389–446.
- 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.
- Liang, P., Petrov, S., Jordan, M. I., and Klein, D. (2007). The infinite pcfg using hierarchical dirichlet processes. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 688–697.
- Liang, P. and Potts, C. (2015). Bringing machine learning and compositional semantics together. *Annual Review of Linguistics*, 1(1):355–376.
- Lieber, R. (2015). *Introducing morphology*. Cambridge University Press.
- Lin, D. (1998). Automatic retrieval and clustering of similar words. In *Proceedings of the 17th international conference on Computational linguistics-Volume 2*, pages 768–774. Association for Computational Linguistics.
- Ling, W., Dyer, C., Black, A., and Trancoso, I. (2015a). Two/too simple adaptations of word2vec for syntax problems. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Ling, W., Luís, T., Marujo, L., Astudillo, R. F., Amir, S., Dyer, C., Black, A. W., and Trancoso, I. (2015b). Finding function in form: Compositional character models for open vocabulary word representation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Ling, X., Singh, S., and Weld, D. S. (2015c). Design challenges for entity linking. *Transactions of the Association for Computational Linguistics*, 3:315–328.
- Linguistic Data Consortium (2005). ACE (automatic content extraction) English annotation guidelines for relations. Technical Report Version 5.8.3, Linguistic Data Consortium.
- 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.
- Loveland, D. W. (2016). *Automated Theorem Proving: a logical basis*. Elsevier.
- Luong, M.-T., Socher, R., and Manning, C. D. (2013). Better word representations with recursive neural networks for morphology. *CoNLL-2013*, 104.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Maas, A. L., Hannun, A. Y., and Ng, A. Y. (2013). Rectifier nonlinearities improve neural network acoustic models. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Mani, I., Verhagen, M., Wellner, B., Lee, C. M., and Pustejovsky, J. (2006). Machine learning of temporal relations. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 753–760.
- Manning, C. D. (2016). Computational linguistics and deep learning. *Computational Linguistics*, 41(4).
- Manning, C. D., Raghavan, P., Schütze, H., et al. (2008). *Introduction to information retrieval*, volume 1. Cambridge university press.
- 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.
- Maron, O. and Lozano-Pérez, T. (1998). A framework for multiple-instance learning. In *Neural Information Processing Systems (NIPS)*, pages 570–576.
- Martins, A. F. T., Smith, N. A., and Xing, E. P. (2009). Concise integer linear programming formulations for dependency parsing. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 342–350.
- Martins, A. F. T., Smith, N. A., Xing, E. P., Aguiar, P. M. Q., and Figueiredo, M. A. T. (2010). Turbo parsers: Dependency parsing by approximate variational inference. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 34–44.
- Martschat, S. and Strube, M. (2015). Latent structures for coreference resolution. *Transactions of the Association for Computational Linguistics*, 3:405–418.
- Matthiessen, C. and Bateman, J. A. (1991). *Text generation and systemic-functional linguistics: experiences from English and Japanese*. Pinter Publishers.
- May, J. (2016). Semeval-2016 task 8: Meaning representation parsing. In *Proceedings of SemEval*, pages 1063–1073.
- McCallum, A. and Li, W. (2003). Early results for named entity recognition with conditional random fields, feature induction and web-enhanced lexicons. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 188–191.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- McCallum, A. and Wellner, B. (2004). Conditional models of identity uncertainty with application to noun coreference. In *NIPS*, pages 905–912.
- McDonald, R., Crammer, K., and Pereira, F. (2005a). Online large-margin training of dependency parsers. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 91–98.
- McDonald, R., Hannan, K., Neylon, T., Wells, M., and Reynar, J. (2007). Structured models for fine-to-coarse sentiment analysis. In *Proceedings of ACL*.
- McDonald, R., Pereira, F., Ribarov, K., and Hajič, J. (2005b). Non-projective dependency parsing using spanning tree algorithms. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 523–530.
- McNamee, P. and Dang, H. T. (2009). Overview of the tac 2009 knowledge base population track. In *Text Analysis Conference (TAC)*, volume 17, pages 111–113.
- Medlock, B. and Briscoe, T. (2007). Weakly supervised learning for hedge classification in scientific literature. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 992–999.
- Miao, Y., Yu, L., and Blunsom, P. (2016). Neural variational inference for text processing. In *Proceedings of the International Conference on Machine Learning (ICML)*.
- Mihalcea, R., Chklovski, T. A., and Kilgarriff, A. (2004). The senseval-3 english lexical sample task. In *Proceedings of SENSEVAL-3*, pages 25–28, Barcelona, Spain. Association for Computational Linguistics.
- Mihalcea, R. and Radev, D. (2011). *Graph-based natural language processing and information retrieval*. Cambridge University Press.
- 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 *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.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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.
- 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)*.
- 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.
- Milne, D. and Witten, I. H. (2008). Learning to link with wikipedia. In *Proceedings of the International Conference on Information and Knowledge Management (CIKM)*, pages 509–518. ACM.
- Minka, T. P. (1999). From hidden markov models to linear dynamical systems. Tech. Rep. 531, Vision and Modeling Group of Media Lab, MIT.
- Minsky, M. (1974). A framework for representing knowledge. Technical Report 306, MIT AI Laboratory.
- Minsky, M. and Papert, S. (1969). *Perceptrons*. MIT press.
- Mintz, M., Bills, S., Snow, R., and Jurafsky, D. (2009). Distant supervision for relation extraction without labeled data. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1003–1011.
- Mirza, P., Sprugnoli, R., Tonelli, S., and Speranza, M. (2014). Annotating causality in the tempeval-3 corpus. In *Proceedings of the EACL 2014 Workshop on Computational Approaches to Causality in Language (CAtoCL)*, pages 10–19.
- Misra, D. K. and Artzi, Y. (2016). Neural shift-reduce ccg semantic parsing. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Mitchell, J. and Lapata, M. (2010). Composition in distributional models of semantics. *Cognitive Science*, 34(8):1388–1429.
- Miwa, M. and Bansal, M. (2016). End-to-end relation extraction using lstms on sequences and tree structures. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1105–1116.
- Mnih, A. and Hinton, G. (2007). Three new graphical models for statistical language modelling. In *Proceedings of the 24th international conference on Machine learning, ICML '07*, pages 641–648, New York, NY, USA. ACM.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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)*.
- Mohammad, S. M. and Turney, P. D. (2013). Crowdsourcing a word–emotion association lexicon. *Computational Intelligence*, 29(3):436–465.
- 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.
- Montague, R. (1973). The proper treatment of quantification in ordinary english. In *Approaches to natural language*, pages 221–242. Springer.
- Morante, R. and Blanco, E. (2012). *sem 2012 shared task: Resolving the scope and focus of negation. In *Proceedings of the First Joint Conference on Lexical and Computational Semantics-Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation*, pages 265–274. Association for Computational Linguistics.
- Morante, R. and Daelemans, W. (2009). Learning the scope of hedge cues in biomedical texts. In *Proceedings of the Workshop on Current Trends in Biomedical Natural Language Processing*, pages 28–36. Association for Computational Linguistics.
- Morante, R. and Sporleder, C. (2012). Modality and negation: An introduction to the special issue. *Computational linguistics*, 38(2):223–260.
- Mostafazadeh, N., Grealish, A., Chambers, N., Allen, J., and Vanderwende, L. (2016). Caters: Causal and temporal relation scheme for semantic annotation of event structures. In *Proceedings of the Fourth Workshop on Events*, pages 51–61, San Diego, California. Association for Computational Linguistics.
- Mueller, T., Schmid, H., and Schütze, H. (2013). Efficient higher-order CRFs for morphological tagging. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 322–332.
- Müller, C. and Strube, M. (2006). Multi-level annotation of linguistic data with mmax2. *Corpus technology and language pedagogy: New resources, new tools, new methods*, 3:197–214.
- (c) Jacob Eisenstein 2014-2017. Work in progress.

- 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.
- 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.*
- Neubig, G., Dyer, C., Goldberg, Y., Matthews, A., Ammar, W., Anastasopoulos, A., Balles-teros, M., Chiang, D., Clothiaux, D., Cohn, T., Duh, K., Faruqui, M., Gan, C., Garrette, D., Ji, Y., Kong, L., Kuncoro, A., Kumar, G., Malaviya, C., Michel, P., Oda, Y., Richardson, M., Saphra, N., Swayamdipta, S., and Yin, P. (2017a). Dynet: The dynamic neural network toolkit.
- Neubig, G., Goldberg, Y., and Dyer, C. (2017b). On-the-fly operation batching in dynamic computation graphs. In *Neural Information Processing Systems (NIPS)*.
- Neuhaus, P. and Bröker, N. (1997). The complexity of recognition of linguistically adequate dependency grammars. In *eacl*, pages 337–343.
- Newman, D., Chemudugunta, C., and Smyth, P. (2006). Statistical entity-topic models. In *Proceedings of Knowledge Discovery and Data Mining (KDD)*, pages 680–686.
- Nguyen, D. and Dogruöz, A. S. (2013). Word level language identification in online multilingual communication. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- 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.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Nivre, J. (2008). Algorithms for deterministic incremental dependency parsing. *Computational Linguistics*, 34(4):513–553.
- Nivre, J., de Marneffe, M.-C., Ginter, F., Goldberg, Y., Hajič, J., Manning, C. D., McDonald, R., Petrov, S., Pyysalo, S., Silveira, N., Tsarfaty, R., and Zeman, D. (2016). Universal dependencies v1: A multilingual treebank collection. In Chair), N. C. C., Choukri, K., Declerck, T., Goggi, S., Grobelnik, M., Maegaard, B., Mariani, J., Mazo, H., Moreno, A., Odijk, J., and Piperidis, S., editors, *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC 2016)*, Paris, France. European Language Resources Association (ELRA).
- Nivre, J., Hall, J., Nilsson, J., Chanev, A., Eryigit, G., Kübler, S., Marinov, S., and Marsi, E. (2007). Maltparser: A language-independent system for data-driven dependency parsing. *Natural Language Engineering*, 13(2):95–135.
- 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.
- Novikoff, A. B. (1962). On convergence proofs on perceptrons. In *Proceedings of the Symposium on the Mathematical Theory of Automata*, volume 12, pages 615–622.
- O'Connor, B., Krieger, M., and Ahn, D. (2010). Tweetmotif: Exploratory search and topic summarization for twitter. In *Proceedings of the International Conference on Web and Social Media (ICWSM)*, pages 384–385.
- Oflazer, K. and Kuruöz, İ. (1994). Tagging and morphological disambiguation of turkish text. In *Proceedings of the fourth conference on Applied natural language processing*, pages 144–149. Association for Computational Linguistics.
- Ohta, T., Tateisi, Y., and Kim, J.-D. (2002). The genia corpus: An annotated research abstract corpus in molecular biology domain. In *Proceedings of the second international conference on Human Language Technology Research*, pages 82–86. Morgan Kaufmann Publishers Inc.
- Onishi, T., Wang, H., Bansal, M., Gimpel, K., and McAllester, D. (2016). Who did what: A large-scale person-centered cloze dataset. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 2230–2235.
- Owoputi, O., O'Connor, B., Dyer, C., Gimpel, K., Schneider, N., and Smith, N. A. (2013). Improved part-of-speech tagging for online conversational text with word clusters. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 380–390.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Packard, W., Bender, E. M., Read, J., Oepen, S., and Drīdan, R. (2014). Simple negation scope resolution through deep parsing: A semantic solution to a semantic problem. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 69–78.
- Paice, C. D. (1990). Another stemmer. In *ACM SIGIR Forum*, volume 24, pages 56–61.
- Pak, A. and Paroubek, P. (2010). Twitter as a corpus for sentiment analysis and opinion mining. In *LREC*, volume 10, pages 1320–1326.
- Palmer, F. R. (2001). *Mood and modality*. Cambridge University Press.
- Palmer, M., Gildea, D., and Kingsbury, P. (2005). The proposition bank: An annotated corpus of semantic roles. *Computational linguistics*, 31(1):71–106.
- Pan, S. J. and Yang, Q. (2010). A survey on transfer learning. *IEEE Transactions on knowledge and data engineering*, 22(10):1345–1359.
- Pan, X., Cassidy, T., Hermjakob, U., Ji, H., and Knight, K. (2015). Unsupervised entity linking with abstract meaning representation. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 1130–1139.
- Pang, B. and Lee, L. (2004). A sentimental education: Sentiment analysis using subjectivity summarization based on minimum cuts. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 271–278.
- 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.
- Parsons, T. (1990). *Events in the Semantics of English*, volume 5. MIT Press.
- Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on Machine Learning (ICML-13)*, pages 1310–1318.
- Peng, F., Feng, F., and McCallum, A. (2004). Chinese segmentation and new word detection using conditional random fields. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, page 562.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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.
- Pereira, F. and Schabes, Y. (1992). Inside-outside reestimation from partially bracketed corpora. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 128–135.
- Pereira, F. C. N. and Shieber, S. M. (2002). *Prolog and natural-language analysis*. Microtome Publishing.
- Peterson, W. W., Birdsall, T. G., and Fox, W. C. (1954). The theory of signal detectability. *Transactions of the IRE professional group on information theory*, 4(4):171–212.
- 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 McDonald, R. (2012). Overview of the 2012 shared task on parsing the web. In *Notes of the First Workshop on Syntactic Analysis of Non-Canonical Language (SANCL)*, volume 59.
- Pinter, Y., Guthrie, R., and Eisenstein, J. (2017). Mimicking word embeddings using subword RNNs. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Plank, B., Søgaard, A., and Goldberg, Y. (2016). Multilingual part-of-speech tagging with bidirectional long short-term memory models and auxiliary loss. In *Proceedings of the Association for Computational Linguistics (ACL)*.
- Polanyi, L. and Zaenen, A. (2006). Contextual valence shifters. In *Computing attitude and affect in text: Theory and applications*. Springer.
- Poon, H. and Domingos, P. (2009). Unsupervised semantic parsing. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1–10.
- Popel, M., Marecek, D., Stepánek, J., Zeman, D., and Zabokrtský, Z. (2013). Coordination structures in dependency treebanks. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 517–527.
- Popescu, A.-M., Etzioni, O., and Kautz, H. (2003). Towards a theory of natural language interfaces to databases. In *Proceedings of Intelligent User Interfaces (IUI)*, pages 149–157.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Poplack, S. (1980). Sometimes i'll start a sentence in spanish y termino en español: toward a typology of code-switching¹. *Linguistics*, 18(7-8):581–618.
- Porter, M. F. (1980). An algorithm for suffix stripping. *Program*, 14(3):130–137.
- Prabhakaran, V., Rambow, O., and Diab, M. (2010). Automatic committed belief tagging. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 1014–1022.
- Pradhan, S., Ramshaw, L., Marcus, M., Palmer, M., Weischedel, R., and Xue, N. (2011). Conll-2011 shared task: Modeling unrestricted coreference in ontonotes. In *Proceedings of the Fifteenth Conference on Computational Natural Language Learning: Shared Task*, pages 1–27. Association for Computational Linguistics.
- Pradhan, S., Ward, W., Hacıoglu, K., Martin, J. H., and Jurafsky, D. (2005). Semantic role labeling using different syntactic views. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 581–588.
- Punyakanok, V., Roth, D., and Yih, W.-t. (2008). The importance of syntactic parsing and inference in semantic role labeling. *Computational Linguistics*, 34(2):257–287.
- Pustejovsky, J., Hanks, P., Sauri, R., See, A., Gaizauskas, R., Setzer, A., Radev, D., Sundheim, B., Day, D., Ferro, L., et al. (2003). The timebank corpus. In *Corpus linguistics*, volume 2003, page 40. Lancaster, UK.
- Pustejovsky, J., Ingria, B., Sauri, R., Castano, J., Littman, J., Gaizauskas, R., Setzer, A., Katz, G., and Mani, I. (2005). The specification language timeml. In *The language of time: A reader*, pages 545–557. Oxford University Press.
- Qiu, G., Liu, B., Bu, J., and Chen, C. (2011). Opinion word expansion and target extraction through double propagation. *Computational linguistics*, 37(1):9–27.
- Quattoni, A., Wang, S., Morency, L.-P., Collins, M., and Darrell, T. (2007). Hidden conditional random fields. *IEEE transactions on pattern analysis and machine intelligence*, 29(10).
- Rabiner, L. (1989). A tutorial on hidden markov models and selected applications in speech recognition. *Proceedings of the IEEE*, 77(2):257–286.
- Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. (2016). Squad: 100,000+ questions for machine comprehension of text. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 2383–2392.
- Rao, D., Yarowsky, D., Shreevats, A., and Gupta, M. (2010). Classifying latent user attributes in twitter. In *Proceedings of Workshop on Search and mining user-generated contents*.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Ratinov, L. and Roth, D. (2009). Design challenges and misconceptions in named entity recognition. In *Proceedings of the Thirteenth Conference on Computational Natural Language Learning*, pages 147–155. Association for Computational Linguistics.
- Ratinov, L., Roth, D., Downey, D., and Anderson, M. (2011). Local and global algorithms for disambiguation to wikipedia. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1375–1384.
- Ratliff, N. D., Bagnell, J. A., and Zinkevich, M. (2007). (approximate) subgradient methods for structured prediction. In *Proceedings of Artificial Intelligence and Statistics (AISTATS)*, pages 380–387.
- Ratnaparkhi, A. (1996). A maximum entropy model for part-of-speech tagging. In *emnlp*, pages 133–142.
- 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.
- Reisinger, D., Rudinger, R., Ferraro, F., Harman, C., Rawlins, K., and Durme, B. V. (2015). Semantic proto-roles. *Transactions of the Association for Computational Linguistics*, 3:475–488.
- 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.
- Ribeiro, F. N., Araújo, M., Gonçalves, P., Gonçalves, M. A., and Benevenuto, F. (2016). Sentibench-a benchmark comparison of state-of-the-practice sentiment analysis methods. *EPJ Data Science*, 5(1):1–29.
- Richardson, M., Burges, C. J., and Renshaw, E. (2013). MCTest: A challenge dataset for the open-domain machine comprehension of text. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 193–203.
- Riedel, S., Yao, L., and McCallum, A. (2010). Modeling relations and their mentions without labeled text. In *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML)*, pages 148–163.
- Riloff, E. (1996). Automatically generating extraction patterns from untagged text. In *Proceedings of the national conference on artificial intelligence*, pages 1044–1049.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Riloff, E. and Wiebe, J. (2003). Learning extraction patterns for subjective expressions. In *Proceedings of the 2003 conference on Empirical methods in natural language processing*, pages 105–112. Association for Computational Linguistics.
- 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.
- Rush, A. M. and Petrov, S. (2012). Vine pruning for efficient multi-pass dependency parsing. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 498–507.
- Rush, A. M., Sontag, D., Collins, M., and Jaakkola, T. (2010). On dual decomposition and linear programming relaxations for natural language processing. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1–11.
- Salinas, E. and Abbott, L. (1996). A model of multiplicative neural responses in parietal cortex. *Proceedings of the national academy of sciences*, 93(21):11956–11961.
- Santos, C. D. and Zadrozny, B. (2014). Learning character-level representations for part-of-speech tagging. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1818–1826.
- Sato, M.-A. and Ishii, S. (2000). On-line em algorithm for the normalized gaussian network. *Neural computation*, 12(2):407–432.
- Saurí, R. and Pustejovsky, J. (2009). Factbank: a corpus annotated with event factuality. *Language resources and evaluation*, 43(3):227.
- Saxe, A. M., McClelland, J. L., and Ganguli, S. (2014). Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Schank, R. C. and Abelson, R. (1977). *Scripts, goals, plans, and understanding*. Erlbaum, Hillsdale, NJ.
- Schapire, R. E. and Singer, Y. (2000). Boostexter: A boosting-based system for text categorization. *Machine learning*, 39(2-3):135–168.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Schnabel, T., Labutov, I., Mimno, D., and Joachims, T. (2015). Evaluation methods for unsupervised word embeddings. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 298–307.
- Schneider, N., Flanigan, J., and O’Gorman, T. (2015). The logic of amr: Practical, unified, graph-based sentence semantics for nlp. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 4–5.
- Schütze, H. (1998). Automatic word sense discrimination. *Computational linguistics*, 24(1):97–123.
- Sennrich, R., Haddow, B., and Birch, A. (2016). Neural machine translation of rare words with subword units. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1715–1725.
- Settles, B. (2012). Active learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 6(1):1–114.
- 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.
- Shen, W., Wang, J., and Han, J. (2015). Entity linking with a knowledge base: Issues, techniques, and solutions. *IEEE Transactions on Knowledge and Data Engineering*, 27(2):443–460.
- 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.
- Smith, D. A. and Eisner, J. (2008). Dependency parsing by belief propagation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 145–156.
- Smith, D. A. and Smith, N. A. (2007). Probabilistic models of nonprojective dependency trees. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 132–140.
- Smith, N. A. (2011). Linguistic structure prediction. *Synthesis Lectures on Human Language Technologies*, 4(2):1–274.
- Snow, R., O’Connor, B., Jurafsky, D., and Ng, A. Y. (2008). Cheap and fast—but is it good?: evaluating non-expert annotations for natural language tasks. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 254–263.
- Socher, R., Bauer, J., Manning, C. D., and Ng, A. Y. (2013a). Parsing with compositional vector grammars. In *Proceedings of the Association for Computational Linguistics (ACL)*.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Socher, R., Huval, B., Manning, C. D., and Ng, A. Y. (2012). Semantic compositionality through recursive matrix-vector spaces. In *Proceedings of the 2012 Joint Conference on Empirical Methods in Natural Language Processing and Computational Natural Language Learning*, pages 1201–1211. Association for Computational Linguistics.
- Socher, R., Perelygin, A., Wu, J. Y., Chuang, J., Manning, C. D., Ng, A. Y., and Potts, C. (2013b). Recursive deep models for semantic compositionality over a sentiment tree-bank. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Søgaard, A. (2013). Semi-supervised learning and domain adaptation in natural language processing. *Synthesis Lectures on Human Language Technologies*, 6(2):1–103.
- Solorio, T. and Liu, Y. (2008). Learning to predict code-switching points. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 973–981. Association for Computational Linguistics.
- Somasundaran, S. and Wiebe, J. (2009). Recognizing stances in online debates. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 226–234.
- 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.
- Sowa, J. F. (2000). *Knowledge representation: logical, philosophical, and computational foundations*. Brooks/Cole, Pacific Grove, CA.
- Spärck Jones, K. (1972). A statistical interpretation of term specificity and its application in retrieval. *Journal of documentation*, 28(1):11–21.
- Spitkovsky, V. I., Alshaw, H., Jurafsky, D., and Manning, C. D. (2010). Viterbi training improves unsupervised dependency parsing. In *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.
- Sproat, R., Gale, W., Shih, C., and Chang, N. (1996). A stochastic finite-state word-segmentation algorithm for chinese. *Computational linguistics*, 22(3):377–404.
- Sra, S., Nowozin, S., and Wright, S. J. (2012). *Optimization for machine learning*. MIT Press.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *The Journal of Machine Learning Research*, 15(1):1929–1958.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Srivastava, R. K., Greff, K., and Schmidhuber, J. (2015). Training very deep networks. In *Neural Information Processing Systems (NIPS)*, pages 2377–2385.
- Steedman, M. and Baldridge, J. (2011). Combinatory categorial grammar. In *Non-Transformational Syntax: Formal and Explicit Models of Grammar*. Wiley-Blackwell.
- Stenetorp, P., Pyysalo, S., Topić, G., Ohta, T., Ananiadou, S., and Tsujii, J. (2012). Brat: a web-based tool for nlp-assisted text annotation. In *Proceedings of the European Chapter of the Association for Computational Linguistics (EACL)*, pages 102–107.
- Stolcke, A., Ries, K., Coccaro, N., Shriberg, E., Bates, R., Jurafsky, D., Taylor, P., Martin, R., Van Ess-Dykema, C., and Meteer, M. (2000). Dialogue act modeling for automatic tagging and recognition of conversational speech. *Computational linguistics*, 26(3):339–373.
- Stone, P. J. (1966). *The General Inquirer: A Computer Approach to Content Analysis*. The MIT Press.
- Stoyanov, V., Gilbert, N., Cardie, C., and Riloff, E. (2009). Conundrums in noun phrase coreference resolution: Making sense of the state-of-the-art. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 656–664.
- Strubell, E., Verga, P., Belanger, D., and McCallum, A. (2017). Fast and accurate entity recognition with iterated dilated convolutions. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Suchanek, F. M., Kasneci, G., and Weikum, G. (2007). Yago: a core of semantic knowledge. In *Proceedings of the Conference on World-Wide Web (WWW)*, pages 697–706.
- Sun, X., Matsuzaki, T., Okanohara, D., and Tsujii, J. (2009). Latent variable perceptron algorithm for structured classification. In *Proceedings of the International Joint Conference on Artificial Intelligence (IJCAI)*, volume 9, pages 1236–1242.
- Sun, Y., Lin, L., Tang, D., Yang, N., Ji, Z., and Wang, X. (2015). Modeling mention, context and entity with neural networks for entity disambiguation. In *IJCAI*, pages 1333–1339.
- Sundermeyer, M., Schlüter, R., and Ney, H. (2012). Lstm neural networks for language modeling. In *INTERSPEECH*.
- Surdeanu, M., Tibshirani, J., Nallapati, R., and Manning, C. D. (2012). Multi-instance multi-label learning for relation extraction. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 455–465.
- Taboada, M., Brooke, J., Tofiloski, M., Voll, K., and Stede, M. (2011). Lexicon-based methods for sentiment analysis. *Computational linguistics*, 37(2):267–307.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Täckström, O., Ganchev, K., and Das, D. (2015). Efficient inference and structured learning for semantic role labeling. *Transactions of the Association for Computational Linguistics*, 3:29–41.
- Täckström, O., McDonald, R., and Uszkoreit, J. (2012). Cross-lingual word clusters for direct transfer of linguistic structure. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 477–487.
- Tang, D., Qin, B., and Liu, T. (2015). Document modeling with gated recurrent neural network for sentiment classification. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1422–1432.
- Tarjan, R. E. (1977). Finding optimum branchings. *Networks*, 7(1):25–35.
- Taskar, B., Guestrin, C., and Koller, D. (2003). Max-margin markov networks. In *Neural Information Processing Systems (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.
- Tesnière, L. (1966). *Éléments de syntaxe structurale*. Klincksieck, Paris, second edition.
- 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.
- Toutanova, K., Klein, D., Manning, C. D., and Singer, Y. (2003). Feature-rich part-of-speech tagging with a cyclic dependency network. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Tromble, R. W. and Eisner, J. (2006). A fast finite-state relaxation method for enforcing global constraints on sequence decoding. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, page 423.
- 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.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Tsvetkov, Y., Faruqui, M., Ling, W., Lample, G., and Dyer, C. (2015). Evaluation of word vector representations by subspace alignment. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 2049–2054.
- Turian, J., Ratinov, L., and Bengio, Y. (2010). Word representations: a simple and general method for semi-supervised learning. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 384–394.
- Turney, P. D. and Pantel, P. (2010). From frequency to meaning: Vector space models of semantics. *Journal of Artificial Intelligence Research*, 37:141–188.
- Twain, M. (1997). *A Tramp Abroad*. Penguin, New York.
- Tzeng, E., Hoffman, J., Darrell, T., and Saenko, K. (2015). Simultaneous deep transfer across domains and tasks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 4068–4076.
- Usunier, N., Buffoni, D., and Gallinari, P. (2009). Ranking with ordered weighted pairwise classification. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1057–1064.
- Uzuner, Ö., Zhang, X., and Sibanda, T. (2009). Machine learning and rule-based approaches to assertion classification. *Journal of the American Medical Informatics Association*, 16(1):109–115.
- Vadas, D. and Curran, J. R. (2011). Parsing noun phrases in the penn treebank. *Computational Linguistics*, 37(4):753–809.
- Van Eynde, F. (2006). NP-internal agreement and the structure of the noun phrase. *Journal of Linguistics*, 42(1):139–186.
- 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.
- Velldal, E., Øvrelid, L., Read, J., and Oepen, S. (2012). Speculation and negation: Rules, rankers, and the role of syntax. *Computational linguistics*, 38(2):369–410.
- Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y., and Manzagol, P.-A. (2010). Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11(Dec):3371–3408.
- Vincze, V., Szarvas, G., Farkas, R., Móra, G., and Csirik, J. (2008). The bioscope corpus: biomedical texts annotated for uncertainty, negation and their scopes. *BMC bioinformatics*, 9(11):S9.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Viterbi, A. (1967). Error bounds for convolutional codes and an asymptotically optimum decoding algorithm. *IEEE transactions on Information Theory*, 13(2):260–269.
- Wager, S., Wang, S., and Liang, P. S. (2013). Dropout training as adaptive regularization. In *Neural Information Processing Systems (NIPS)*, pages 351–359.
- Wainwright, M. J. and Jordan, M. I. (2008). Graphical models, exponential families, and variational inference. *Foundations and Trends® in Machine Learning*, 1(1-2):1–305.
- Wang, C., Xue, N., and Pradhan, S. (2015). A Transition-based Algorithm for AMR Parsing. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 366–375.
- Wang, H., Onishi, T., Gimpel, K., and McAllester, D. (2017). Emergent predication structure in hidden state vectors of neural readers. In *Proceedings of the 2nd Workshop on Representation Learning for NLP*, pages 26–36.
- Weaver, W. (1955). Translation. *Machine translation of languages*, 14:15–23.
- 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.
- Weinberger, K., Dasgupta, A., Langford, J., Smola, A., and Attenberg, J. (2009). Feature hashing for large scale multitask learning. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1113–1120.
- Weston, J., Bengio, S., and Usunier, N. (2011). Wsabie: Scaling up to large vocabulary image annotation. In *IJCAI*, volume 11, pages 2764–2770.
- Wiebe, J., Wilson, T., and Cardie, C. (2005). Annotating expressions of opinions and emotions in language. *Language resources and evaluation*, 39(2):165–210.
- Wieting, J., Bansal, M., Gimpel, K., and Livescu, K. (2015). Towards universal paraphrastic sentence embeddings. *arXiv preprint arXiv:1511.08198*.
- Wieting, J., Bansal, M., Gimpel, K., and Livescu, K. (2016). CHARAGRAM: Embedding words and sentences via character n-grams. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1504–1515.
- 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.
- Winograd, T. (1972). Understanding natural language. *Cognitive psychology*, 3(1):1–191.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Wiseman, S., Rush, A. M., and Shieber, S. M. (2016). Learning global features for coreference resolution. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 994–1004.
- Wolfe, T., Dredze, M., and Van Durme, B. (2017). Pocket knowledge base population. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 305–310.
- Wong, Y. W. and Mooney, R. J. (2006). Learning for semantic parsing with statistical machine translation. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*, pages 439–446.
- Wu, B. Y. and Chao, K.-M. (2004). *Spanning trees and optimization problems*. CRC Press.
- Wu, F. and Weld, D. S. (2010). Open information extraction using wikipedia. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 118–127.
- Xia, F. (2000). The part-of-speech tagging guidelines for the penn chinese treebank (3.0). Technical report, University of Pennsylvania Institute for Research in Cognitive Science.
- Xu, W., Liu, X., and Gong, Y. (2003). Document clustering based on non-negative matrix factorization. In *SIGIR*, pages 267–273. ACM.
- Xu, Y., Mou, L., Li, G., Chen, Y., Peng, H., and Jin, Z. (2015). Classifying relations via long short term memory networks along shortest dependency paths. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1785–1794.
- Xue, N. et al. (2003). Chinese word segmentation as character tagging. *Computational Linguistics and Chinese Language Processing*, 8(1):29–48.
- Yamada, H. and Matsumoto, Y. (2003). Statistical dependency analysis with support vector machines. In *Proceedings of IWPT*, volume 3, pages 195–206.
- Yang, Y., Chang, M.-W., and Eisenstein, J. (2016). Toward socially-infused information extraction: Embedding authors, mentions, and entities. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Yang, Y. and Eisenstein, J. (2013). A log-linear model for unsupervised text normalization. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*.
- Yang, Y. and Eisenstein, J. (2015). Unsupervised multi-domain adaptation with feature embeddings. In *Proceedings of the North American Chapter of the Association for Computational Linguistics (NAACL)*.
- Yang, Y., Yih, W.-t., and Meek, C. (2015). WikiQA: A challenge dataset for open-domain question answering. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 2013–2018.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- 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.
- Yu, C.-N. J. and Joachims, T. (2009). Learning structural svms with latent variables. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 1169–1176.
- Yu, F. and Koltun, V. (2016). Multi-scale context aggregation by dilated convolutions. In *Proceedings of the International Conference on Learning Representations (ICLR)*.
- Zeiler, M. D. (2012). Adadelta: an adaptive learning rate method. *arXiv preprint arXiv:1212.5701*.
- Zelenko, D., Aone, C., and Richardella, A. (2003). Kernel methods for relation extraction. *The Journal of Machine Learning Research*, 3:1083–1106.
- Zelle, J. M. and Mooney, R. J. (1996). Learning to parse database queries using inductive logic programming. In *Proceedings of the National Conference on Artificial Intelligence (AAAI)*, pages 1050–1055.
- Zeng, D., Liu, K., Lai, S., Zhou, G., and Zhao, J. (2014). Relation classification via convolutional deep neural network. In *Proceedings of the International Conference on Computational Linguistics (COLING)*, pages 2335–2344.
- Zettlemoyer, L. S. and Collins, M. (2005). Learning to map sentences to logical form: Structured classification with probabilistic categorical grammars. In *Proceedings of UAI*.
- 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.
- Zhang, X., Zhao, J., and LeCun, Y. (2015). Character-level convolutional networks for text classification. In *Neural Information Processing Systems (NIPS)*, pages 649–657.
- Zhang, Y. and Clark, S. (2008). A tale of two parsers: investigating and combining graph-based and transition-based dependency parsing using beam-search. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 562–571.
- Zhang, Y., Lei, T., Barzilay, R., Jaakkola, T., and Globerson, A. (2014). Steps to excellence: Simple inference with refined scoring of dependency trees. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 197–207.
- Zhang, Y. and Nivre, J. (2011). Transition-based dependency parsing with rich non-local features. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 188–193.

(c) Jacob Eisenstein 2014-2017. Work in progress.

- Zhou, J. and Xu, W. (2015). End-to-end learning of semantic role labeling using recurrent neural networks. In *Proceedings of the Association for Computational Linguistics (ACL)*, pages 1127–1137.
- Zhu, J., Nie, Z., Liu, X., Zhang, B., and Wen, J.-R. (2009). Statsnowball: a statistical approach to extracting entity relationships. In *Proceedings of the Conference on World-Wide Web (WWW)*, pages 101–110.
- Zhu, X., Ghahramani, Z., and Lafferty, J. D. (2003). Semi-supervised learning using gaussian fields and harmonic functions. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 912–919.
- Zhu, X. and Goldberg, A. B. (2009). Introduction to semi-supervised learning. *Synthesis lectures on artificial intelligence and machine learning*, 3(1):1–130.
- Zou, W. Y., Socher, R., Cer, D., and Manning, C. D. (2013). Bilingual word embeddings for phrase-based machine translation. In *Proceedings of Empirical Methods for Natural Language Processing (EMNLP)*, pages 1393–1398.

Index

- F*-measure, 74
- K*-means, 89
- α -conversion, 271
- β -conversion, 269
- β -reduction, 269
- n*-gram language models, 188
- n*-gram, 21

- ablation tests, 76
- Abstract Meaning Representation (AMR), 285, 299
- accepting path, 183
- access (an entity), 336
- accuracy, 20, 73
- activation function, 125
- Active learning, 109
- AdaDelta, 55
- AdaGrad, 35, 55
- Adam, 55
- adjectives, 169
- adjuncts (semantics), 284, 288
- adpositions, 169
- adverbs, 169
- adversarial networks, 105
- affix (morphology), 186
- agent (thematic role), 286
- alignment, 301
- Amazon Mechanical Turk, 82
- ambiguity, 208
- ambiguity (grammar), 200
- anaphora resolution, 335

- antecedent (coreference), 344
- antecedent structure, 344
- antonymy, 66
- apophony, 184
- Arc-eager dependency parsing, 256
- arc-factored, 247
- arc-standard, 254
- area under the curve (AUC), 75
- arguments, 353
- article (syntax), 174
- aspect, 169
- Attachment ambiguity, 220
- attention mechanism, 376
- autoencoders, 329
- automated theorem provers, 264
- automatic differentiation, 50
- auxiliary verbs, 170
- average mutual information, 317
- averaged perceptron, 23

- backchannel, 179
- backpropagation, 49, 126
- backpropagation through time, 128
- bag of words, 11
- balanced test set, 73
- batch normalization, 54
- batch optimization, 34
- Baum-Welch algorithm, 162
- Bayes' rule, 387
- Bayesian nonparametrics, 95
- Beam search, 256

- best-path algorithm, 188
- bias, 117
- bias (learning theory), 20
- bias-variance, 119
- bias-variance tradeoff, 20
- bidirectional recurrent neural network, 161
- bigrams, 21, 62
- bilexical, 234
- Bilexical features, 247
- bilinear product, 314
- binarization (context-free grammar), 201
- binarize, 219
- binding (anaphora), 337
- binomial distribution, 76
- binomial random variable, 392
- binomial test, 76
- BIO notation, 176, 296
- biomedical natural language processing, 175
- bipartite graph, 303
- bitext, 116
- Bonferroni correction, 79
- boolean semiring, 189
- boosting, 42
- bootstrap samples, 78
- bound variable, 266
- Brown clusters, 311
- byte-pair encodings, 327
- c-command, 337
- case, 173
- case marking, 210
- Catalan number, 217
- cataphora, 336
- center embedding, 198
- centering theory, 336, 339
- chain FSA, 196
- chain rule, 387
- chance agreement, 81
- character-level language models, 133
- chart parsing, 217
- Chomsky Normal Form (CNF), 201
- CKY algorithm, 217
- class imbalance, 73
- classification weights, 11
- cleft, 341
- closed-vocabulary, 132
- closure (regular languages), 182
- cloze question answering, 375
- clustering, 88
- co-training, 100
- code switching, 177
- code-switching, 171
- Cohen's Kappa, 81
- collapsed dependency grammar, 244
- collapsed Gibbs sampling, 108
- collective entity linking, 358
- collocation extraction, 328
- collocation features, 68
- combinatory categorial grammar, 211
- complement clause, 203
- complement event (probability), 386
- Complement structure, 221
- composition (CCG), 212
- compositionality, 325
- computation graph, 42, 49, 128
- concept (AMR), 299
- conditional independence, 390
- conditional log-likelihood, 46
- conditional probability, 31, 387
- conditional probability distribution, 391
- Conditional Random Field, 155
- conditionally independent, 144
- confidence interval, 78
- connected (graph theory), 302
- consistency (logic), 267
- constants (logic), 262
- constituents, 202
- constrained optimization, 28, 293
- Constraint-driven learning, 109
- constraints, 293

(c) Jacob Eisenstein 2014-2017. Work in progress.

- content words, 169
- context-free grammars (CFGs), 199
- context-free languages, 197, 199
- context-sensitive languages, 210
- continuous bag-of-words (CBOW), 319
- contradiction, 330
- conversational turns, 179
- convex, 25
- convex optimization, 34
- convexity, 52, 277
- convolutional neural nets, 162
- convolutional neural network, 133
- convolutional neural networks, 47, 56, 62, 177, 296, 365
- coordinate ascent, 93
- coordinating conjunctions, 170
- Coordination scope, 221
- copula, 207
- copula verb, 169
- Coreference resolution, 335
- coreference resolution, 331
- cosine similarity, 324
- cost-augmented decoding, 30, 154
- cost-augmented inference, 30
- count, 72
- critical point, 52
- cross-document coreference resolution, 356
- cross-entropy, 46, 366
- cross-serial dependencies, 210
- cross-validation, 20
- crowdsourcing, 82
- cumulative probability distribution, 77
- dead neurons, 45
- decidability (logic), 267
- decision trees, 42
- deep learning, 41
- definiteness, 174
- denotation (semantics), 262
- dependency grammar, 242
- dependency graph, 242
- dependency path, 68, 292, 363
- dependent, 242
- derivation (context-free languages), 199
- derivation (semantic parsing), 270, 274
- derivational ambiguity, 213
- derivational morphology, 184
- determiner, 170
- determiner phrase, 205
- deterministic FSA, 183
- development set, 20, 73
- dialogue acts, 82, 178
- digital humanities, 61
- dilated convolution, 58
- dilated convolutions, 177
- Dirichlet distribution, 107
- discount, 120
- discourse relations, 331
- discrete random variable, 390
- discriminative learning, 22
- disjoint events, 386
- Distant supervision, 109
- distant supervision, 368, 369
- distributional, 310
- distributional hypothesis, 309, 310
- distributional statistics, 67
- document frequency, 357
- domain adaptation, 87, 102
- dropout, 51, 128
- dual decomposition, 296
- dynamic computation graphs, 50
- dynamic programming, 139
- dynamic semantics, 279
- E-step, 91
- early stopping, 24, 55
- edit distance, 191
- effective counts, 120
- elementwise nonlinearity, 44
- embedding, 319
- emotion, 64

- empty string, 182
- encoder-decoder model, 329
- ensemble, 299
- ensemble methods, 42
- entailment, 267, 330
- entities, 353
- entity (coreference), 343
- entity embeddings, 358
- Entity linking, 335
- entity linking, 353, 355, 366
- entropy, 38, 91
- estimation, 392
- event, 370
- event (probability), 385
- event coreference, 371
- event detection, 370
- event semantics, 283
- events, 353
- evidentiality, 174, 373
- evoke (an entity), 336
- exchange clustering, 318
- existential quantifier, 266
- expectation, 391
- expectation maximization, 89, 121
- expectation semiring, 197
- explicit semantic analysis, 312
- extra-propositional semantics, 372
- extractive question-answering, 376
- extraposition, 341

- factoid questions, 303
- factoids, 374
- factuality, 373
- false discovery rate, 79
- False negative, 73
- False positive, 73
- false positive, 389
- false positive rate, 75
- feasible set, 293
- feature co-adaptation, 51
- feature function, 11, 21

- feature hashing, 72
- feature noising, 52
- feature selection, 37
- feature templates, 153
- feedforward neural network, 44
- fine-tuned word embeddings, 325
- finite-state acceptor (FSA), 183
- finite-state automata, 183
- finite-state composition, 194
- finite-state transducers, 186, 191
- first-order logic, 264
- fluent, 115
- focus, 301
- formal language theory, 181
- forward algorithm, 229
- forward variables, 148
- forward-backward algorithm, 158, 197
- frame elements, 289
- FrameNet, 289
- frames, 288
- free (anaphora), 337
- free variable, 265
- Frobenius norm, 51
- function (first-order logic), 266
- function words, 169
- functional margin, 28

- garden path sentence, 136
- gate (neural networks), 46
- gazetteer, 363
- gazetteers, 342
- gazzeteers, 176
- generalization, 25
- generative model, 14
- generative process, 122
- geometric margin, 28
- Gibbs Sampling, 106, 164
- Gibbs sampling, 359
- gloss, 116
- gradient, 26
- gradient clipping, 54

(c) Jacob Eisenstein 2014-2017. Work in progress.

- gradient descent, 34
- Gram matrix, 363
- graph-based dependency parsing, 246
- graphical model, 144
- graphics processing units (GPUs), 162, 177
- grid search, 20
- Hamming cost, 154
- hanzi, 69
- hard expectation-maximization, 94
- hard tanh, 297
- head, 242, 363
- head word, 203, 342
- head words, 241
- hedging, 373
- held-out data, 131
- Hessian matrix, 34
- hidden Markov models, 144
- hidden variable perceptron, 197
- hierarchical clustering, 315
- hierarchical softmax, 126, 321
- highway network, 46
- hinge loss, 25
- human computation, 82
- hyperparameter, 19
- hyponymy, 67
- illocutionary force, 178
- incremental expectation maximization, 94
- independent and identically distributed (IID), 14
- Indicator random variables, 390
- inference, 137
- inference (logic), 261
- inference rules, 264
- inflectional affixes, 70
- inflectional morphology, 169, 184, 192
- information extraction, 353
- information retrieval, 367
- input word embeddings, 125
- inside algorithm, 229
- inside-outside algorithm, 237
- instance (AMR), 299
- instance labels, 14
- integer linear program, 378
- Integer linear programming, 346
- integer linear programming, 293, 359
- inter-annotator agreement, 81
- interjections, 169
- interpolated n -gram language model, 189
- interpolation, 121
- interval algebra, 371
- inverse document frequency, 357
- inverse relation (AMR), 301
- inversion (finite-state), 194
- irrealis, 62
- Jeffreys-Perks law, 120
- Jensen's inequality, 91
- joint probabilities, 391
- joint probability, 14, 31
- Kalman Smoother, 163
- Katz backoff, 120
- kernel function, 363
- kernel methods, 42
- kernel support vector machine, 42, 364
- Kleene star, 182
- knowledge base, 353
- knowledge base population, 366
- L-BFGS, 35
- label propagation, 101, 111
- labeled dependencies, 243
- lambda calculus, 268
- lambda expressions, 269
- language (formal language theory), 181
- language model, 116
- Laplace smoothing, 19, 120
- large margin classification, 27

- latent conditional random fields, 277
- latent semantic analysis, 312, 314
- latent variable, 90, 197, 369
- latent variable perceptron, 277
- latent variables, 276
- layer normalization, 55
- leaky ReLU, 45
- least squares, 64
- leave-one-out, 21
- lemma, 65
- lemma (lexical semantics), 192
- lemmatization, 70
- Levenshtein edit distance, 191
- lexical entry, 270
- lexical semantics, 65
- lexical unit (frame semantics), 288
- lexicalized tree-adjoining grammar, 235
- lexicon, 270
- lexicon (CCG), 212
- lexicon-based classification, 64
- lexicon-based sentiment analysis, 62
- Lidstone smoothing, 120
- light verb, 301
- likelihood, 388
- linear regression, 64
- linear separability, 22
- link functions, 37
- literal character, 182
- local optimum, 93
- log-bilinear language model, 326
- logistic function, 37
- logistic loss, 32
- Logistic regression, 31
- logistic regression, 37
- Long short-term memories, 126
- long short-term memories, 298
- long short-term memory, 129
- long short-term memory (LSTM), 46, 173
- lookup layer, 47
- loss function, 25
- LSTM, 129
- LSTM-CRF, 161, 298
- machine reading, 375
- Macro F -measure, 74
- macro-reading, 354
- margin, 22, 27
- marginal probability distribution, 391
- marginalize, 387
- markable, 343
- Markov blanket, 144
- Markov Chain Monte Carlo (MCMC), 96, 106, 164
- matrix-tree theorem, 249
- max pooling, 297
- max-margin Markov network, 153
- max-product algorithm, 147
- maximum a posteriori, 19, 393
- maximum conditional likelihood, 31
- maximum entropy, 37
- maximum likelihood, 392
- maximum likelihood estimate, 18
- maximum likelihood estimation, 14
- McNemar's test, 76
- meaning representation, 261
- membership problem, 181
- mention (coreference), 343
- mention (entity), 353
- mention (information extraction), 355
- mention ranking, 344
- mention-pair model, 344
- meronymy, 67
- micro F -measure, 74
- micro-reading, 354
- mildly context-sensitive, 210
- minibatch, 35
- minimization (FSA), 186
- modality, 372
- model builder, 267
- model checker, 267
- model-theoretic semantics, 262
- modeling (machine learning), 34

- Modifier scope, 221
- modus ponens, 264
- moment-matching, 37
- monomorphemic, 186
- morphemes, 133, 186
- morphological analysis, 193
- morphological generation, 193
- Morphology, 151
- morphology, 71, 184, 325
- morphosyntactic, 168
- morphosyntactic attributes, 172
- morphotactic, 185
- multi-task learning, 297
- multilayer perceptron, 44
- multinomial distribution, 15
- multinomial naïve Bayes, 16
- Multiple instance learning, 109
- multiple instance learning, 368
- Multitask learning, 109

- Naïve Bayes, 15
- name dictionary, 356
- named entities, 175
- named entity linking, 355
- named entity recognition, 161, 353, 355
- named entity types, 355
- narrow convolution, 57
- nearest-neighbor, 42, 364
- negation, 62, 372
- negative sampling, 321, 322, 361
- Neo-Davidsonian event semantics, 284
- neural networks, 42, 124
- NIL entity, 356
- noise-contrastive estimation, 126
- noisy channel model, 116
- nominal modifier, 205
- nominals, 335
- nominals (coreference), 342
- non-convex, 25
- non-core roles (AMR), 300
- non-projective dependency parsing, 249

- non-terminals, 199
- normalization, 70
- noun phrase, 203
- nouns, 168
- NP-hard, 37, 359
- null hypothesis, 76
- numeral (part of speech), 171

- offset feature, 13
- one-dimensional convolution, 56
- one-tailed p-value, 77
- one-versus-all multiclass classification, 364
- one-versus-one multiclass classification, 364
- online expectation maximization, 94
- online learning, 22, 35
- open information extraction, 369
- open word classes, 168
- opinion polarity, 61
- oracle, 303
- orthography, 186, 195
- orthonormal matrix, 53
- out-of-vocabulary words, 172
- overfit, 20
- overfitting, 24
- overgenerate, 272
- overgeneration, 194, 203

- parameters, 392
- paraphrase, 330
- parent annotation, 231
- parsing, 199
- part-of-speech, 167
- part-of-speech tagging, 135
- particle (part-of-speech), 171, 207
- Particle versus preposition, 221
- partition, 387
- partition function, 158
- path (finite-state automata), 183
- Penn Treebank, 202, 227

- perceptron, 22
- perplexity, 131
- phonology, 186
- phrase (syntax), 203
- phrase-structure grammar, 202
- pivot features, 104
- pleonastic, 341
- pointwise mutual information, 314
- polysemous, 66
- pooling (convolution), 57
- pooling (in convolutional neural networks), 298
- positive pointwise mutual information, 315
- posterior, 388
- potentials, 158
- pre-trained word embeddings, 297
- pre-trained word representations, 324
- precision, 74
- precision-recall curve, 368
- predicate, 353
- predicative adjectives, 207
- predictive likelihood, 95
- prepositional phrase, 207
- presence, 72
- principle of compositionality, 268
- prior, 388
- prior expectation, 393
- probabilistic context-free grammar, 210, 224
- probabilistic models, 392
- probabilistic topic model, 359
- probability distribution, 390
- probability mass function, 77
- probability simplex, 15
- processes, 372
- production rules, 199
- productivity, 185
- projection function, 104
- projectivity, 245
- pronoun, 170
- PropBank, 288
- proper nouns, 169
- property (logic), 265
- proposition, 372
- propositions, 262, 263
- prosody, 179
- proto-roles, 287
- pumping lemma, 197
- pushdown automata, 199
- quadratic program, 28
- quasi-Newton optimization, 35
- question answering, 329, 355
- random outcomes, 385
- ranking, 64, 356
- ranking loss, 356
- recall, 74
- receiver operating characteristic (ROC), 75
- rectified linear unit, 125
- rectified linear unit (ReLU), 45, 297
- recurrent neural network, 125, 298
- recurrent neural networks, 365
- recursive neural networks, 327, 330
- reference arguments, 306
- regression, 64
- regular expression, 182
- regular language, 182
- regularization, 30
- reification (events), 283
- relation extraction, 304, 361
- relations, 353
- relations (information extraction), 353
- relations (logic), 262
- relative frequency estimate, 18, 393
- ReLU, 125
- reranking, 239
- residual networks, 46
- retrofitting (word embeddings), 327
- ridge regression, 64

(c) Jacob Eisenstein 2014-2017. Work in progress.

- root (morpheme), 327
- saddle points, 52
- sample space, 385
- satisfaction (logic), 267
- schema, 353, 354, 369
- second-order logic, 265
- seed lexicon, 65
- self-training, 100
- semantic, 168
- semantic concordance, 68
- semantic parsing, 268
- semantic role, 284
- Semantic role labeling, 284
- semantic role labeling, 362, 370
- semantic underspecification, 279
- semantics, 230
- semi-supervised learning, 87, 97, 324
- semiring algebra, 150
- semiring notation, 189
- semisupervised, 68
- senses, 287
- sentence (logic), 266
- sentiment, 61
- sentiment lexicon, 14
- shift-reduce, 217
- shifted positive pointwise mutual information, 322
- shortest-path algorithm, 188
- sigmoid, 43
- sigmoid function, 125
- simplex, 107
- singular value decomposition, 53
- singular value decomposition (SVD), 96
- singular vectors, 54
- skipgram word embeddings, 320
- slack variables, 29
- slot filling, 366
- smooth functions, 26
- smoothing, 119
- soft K -means, 89
- softmax, 43, 124, 366
- source domain, 102
- sparse matrix, 314
- sparsity, 37
- spectral learning, 94
- speech acts, 178
- split constituents, 290
- Spurious ambiguity, 255
- spurious ambiguity, 213, 274
- squashing functions, 45
- stand-off annotations, 81
- Stanford Natural Language Inference corpus, 330
- stationary point, 52
- statistical learning theory, 23
- statistical significance, 76
- stem (morphology), 186
- stemmer, 70
- step size, 34
- stochastic gradient descent, 26, 35
- stoplist, 72
- stopwords, 72
- string (formal language theory), 181
- strongly equivalent grammars, 201
- structured classification, 12
- structured perceptron, 152
- structured support vector machine, 153
- subgradient, 26, 37
- subjectivity detection, 63
- subordinating conjunctions, 170
- sum-product, 149
- supersenses, 324
- supervised machine learning, 14
- support vector machine, 28
- support vectors, 28
- surface form, 194
- synonymy, 66, 309
- synsets, 66
- syntactic dependency, 241
- syntactic path, 291
- syntax, 167, 202

- tanh activation function, 45
- target domain, 102
- Targeted sentiment analysis, 63
- tense, 169
- test set, 20
- test statistic, 76
- text classification, 11
- thematic roles, 286
- third axiom of probability, 386
- TimeML, 371
- token, 115
- tokenization, 69, 177
- Tokenizers, 115
- tokens, 16
- trace (syntax), 213
- training set, 14
- transduction, 182
- Transfer learning, 109
- translation model, 116
- Tree-Adjoining Grammar, 228
- tree-adjoining grammar, 211
- trellis, 140, 196
- trigrams, 21
- tropical semiring, 150, 189
- True negative, 73
- True positive, 73
- true positive, 389
- true positive rate, 75
- truncated singular value decomposition, 314
- truth conditions, 266
- tuning set, 20, 73
- two-tailed test, 77
- type systems, 272
- type-raising, 212, 272
- types, 16
- unary productions, 200
- underfit, 20
- underflow, 14
- undergeneration, 194, 203
- Universal Dependencies, 168
- universal quantifier, 266
- unsupervised, 68
- unsupervised learning, 63, 87
- utterances, 178
- validation function (semantic parsing), 278
- validity (logic), 267
- vanishing gradient, 45
- variable (AMR), 299
- variance, 79
- variance (learning theory), 19
- verb phrase, 203
- VerbNet, 286
- verbs, 169
- vertical Markovization, 231
- Viterbi algorithm, 139
- Viterbi variable, 140
- WARP loss, 360
- weakly equivalent grammars, 201
- weight decay, 51
- weighted context-free grammar, 224
- weighted finite-state acceptors, 187
- wide convolution, 57
- Wikification, 355
- word embedding, 47
- word embeddings, 41, 126, 310, 311
- word representations, 310
- word sense disambiguation, 65
- word sense induction, 87
- word senses, 65
- word tokens, 69
- world model, 262
- yield (context-free grammars), 199
- zero-one loss, 25