

CONTENTS

1	NLP with Deep Learning	2
1.1	Word Vector Representations (Lec 2)	3
1.2	GloVe (Lec 3)	6
2	Speech and Language Processing	9
2.1	Introduction (Ch. 1 2nd Ed.)	10
2.2	Morphology (Ch. 3 2nd Ed.)	11
2.3	N-Grams (Ch. 6 2nd Ed.)	12
2.4	Naive Bayes and Sentiment (Ch. 6 3rd Ed.)	14
2.5	Hidden Markov Models (Ch. 9 3rd Ed.)	16
2.6	POS Tagging (Ch. 10 3rd Ed.)	19
2.7	Formal Grammars (Ch. 11 3rd Ed.)	22
2.8	Vector Semantics (Ch. 15)	23
2.9	Semantics with Dense Vectors (Ch. 16)	26
2.10	Information Extraction (Ch. 21 3rd Ed)	29

NLP WITH DEEP LEARNING

CONTENTS

1.1	Word Vector Representations (Lec 2)	3
1.2	GloVe (Lec 3)	6

Word Vector Representations (Lec 2)

Meaning of a word. Common answer is to use a *taxonomy* like WordNet that has hypernyms (is-a) relationships. Problems with this discrete representation: misses nuances, e.g. the words in a set of synonyms can actually have quite different meanings/connotations. Furthermore, viewing words as atomic symbols is a bit like using one-hot vectors of words in a vocabulary space (inefficient).

Distributed representations. Want a way to encode word vectors such that two similar words have a similar structure/representation. The **distributional similarity-based**¹ approach represents words by means of its *neighbors* in the sentences in which it appears. You end up with a dense “vector for each word type, chosen so that it is good at predicting other words appearing in its context.”

Skip-gram prediction. Given a word w_t at position t in a sentence, learn to predict [probability of] some number of surrounding words, given w_t . Standard minimization with negative log-likelihood:

$$J(\theta) = -\frac{1}{T} \sum_{t=1}^T \sum_{-m \leq j \leq m} \log \Pr(w_{t+j} | w_t) \quad (1)$$

$$\Pr(o | c) = \frac{e^{\mathbf{u}_o^T \mathbf{v}_c}}{\sum_{w=1}^{\text{vocab size}} e^{\mathbf{u}_w^T \mathbf{v}_c}} \quad (2)$$

I cannot believe this actually works.

where

- The params θ are the vector representation of the words (they are the *only* learnable parameters here).
- m is our radius/window size.
- o and c are indices into our vocabulary (somewhat inconsistent notation).
- Yes, they are using different vector representations for \mathbf{u} (context words) and \mathbf{v} (center words). I’m assuming one reason this is done is because it makes the model architecture simpler/easier to build.

Some subtleties:

¹Note that this is distinct from the way “distributed” is meant in “distributed representation.” In contrast, distributional similarity-based representations refers to the notion that you can describe the meaning of words by understanding the context in which they appear.

- Looks like e.g. $Pr(w_{t+j} | w_t)$ doesn't really care what the value of j is, it is just modeling the probability that it is *somewhere* in the context window. The w_t are one-hot vectors into the vocabulary. Standard tricks for simplifying the cross-entropy loss apply.
- Equation 2 should be interpreted as the probability that the o^{th} word in our vocabulary occurs in the context window of the c^{th} word in our vocabulary.
- The model architecture is *identical* to an autoencoder. However, the (big) difference is the training procedure and interpretation of the model parameters going “in” versus the parameters going “out”.

Sentence Embeddings. It turns out that simply taking a weighted average of word vectors and doing some PCA/SVD is a competitive way of getting unsupervised word embeddings. Apparently it *beats* supervised learning with LSTMs (?!). The authors claim the theoretical explanation for this method lies in a latent variable generative model for sentences (of course). Approach:

Discussion based on paper by Arora et al. (2017).

1. Compute the weighted average of the word vectors in the sentence:

$$\frac{1}{N} \sum_i^N \frac{a}{a + p(\mathbf{w}_i)} \mathbf{w}_i \quad (3)$$

The authors call their weighted average the **Smooth Inverse Frequency (SIF)**.

where \mathbf{w}_i is the word vector for the i th word in the sentence, a is a parameter, and $p(\mathbf{w}_i)$ is the (estimated) word frequency [over the entire corpus].

2. Remove the projections of the average vectors on their first principal component (“common component removal”) (y tho?).

1

Algorithm 1 Sentence Embedding

Input: Word embeddings $\{v_w : w \in \mathcal{V}\}$, a set of sentences \mathcal{S} , parameter a and estimated probabilities $\{p(w) : w \in \mathcal{V}\}$ of the words.

Output: Sentence embeddings $\{v_s : s \in \mathcal{S}\}$

- 1: **for all** sentence s in \mathcal{S} **do**
 - 2: $v_s \leftarrow \frac{1}{|s|} \sum_{w \in s} \frac{a}{a + p(w)} v_w$
 - 3: **end for**
 - 4: Compute the first principal component u of $\{v_s : s \in \mathcal{S}\}$
 - 5: **for all** sentence s in \mathcal{S} **do**
 - 6: $v_s \leftarrow v_s - uu^\top v_s$
 - 7: **end for**
-

Further Reading.

- Learning representations by back-propagating errors (Rumelhard et al., 1986)
- A Neural Probabilistic Language Model (Bengio et al., 2003)
- NLP (almost) from Scratch (Collobert & Watson, 2008)
- Word2Vec (Miklov et al. 2013)

GloVe (Lec 3)

Skip-gram and negative sampling. Main idea:

- Split the loss function from last lecture into two (additive) terms corresponding to the numerator and denominator respectively (you’ve done this a trillion times).
- The second term is an expectation over all the words in your vocab space. That is huge, so instead we only use a subsample of size k (the negative samples).
- Interpretation: First term is maximizing $\Pr(o | c)$, the probability that the true outside word (given by index o) occurs given context (index) c . Second term is minimizing the probability of random words (the negative samples) occurring around the center (context) word given by c .

To sample the negative samples, draw from $P(w) = U(w)^{3/4}/Z$, where U is the **unigram distribution**.

GloVe (Global Vectors). Given some co-occurrence matrix we computed with previous methods, we can use the following GloVe loss function over all pairs of co-occurring words in our matrix:

$$J(\theta) = \sum_{i,j=1}^W f(P_{ij})(u_i^T v_j - \log P_{ij})^2 \quad (4)$$

where P_{ij} is computed simply from the counts of words i and j co-occurring (empirical probability) and f is some squashing function that really isn’t discussed in this lecture (**TODO**).

Evaluating word vectors.

- **Word Vector Analogies:** Basically, determining if we can do standard analogy fill-in-the-blank problems: “*man [a] is to woman [b] as king [c] is to <blank>*” (if you answered “queen”, you’d make a good AI). We can determine this using a standard cosine distance measure:

$$d = \arg \max_i \frac{(x_b - x_a + x_c)^T x_i}{\|x_b - x_a + x_c\|} \quad (5)$$

Woah that is pretty neat. The solution is $x_i = \text{queen}$. $x_b - x_a$ is the vector pointing from man to woman, which encodes the type of similarity we are looking for with the other pair. Therefore, we take the vector to “king” and *add* the aforementioned difference vector – the resultant vector should point to “queen”. Neat!

Derivation. Based on the descriptions in the original paper² We want to develop a model for learning word vectors.

1. The authors argue that “the appropriate starting point for word vector learning should be with ratios of co-occurrence probabilities rather than the probabilities themselves.” The most general such model takes the form,

$$F(w_i, w_j, \tilde{w}_k) = \frac{\Pr[\tilde{w}_k | w_i]}{\Pr[\tilde{w}_k | w_j]} \equiv \frac{P_{ik}}{P_{jk}} \quad \text{where all } w \in \mathbb{R}^d \quad (6)$$

and the tilde in \tilde{w}_k denotes that \tilde{w}_k is a *context* word vector, which are given a distinct space from the word vectors w_i and w_j being compared. We compute all P_{ik} via frequency counts over the corpus.

2. Now that we’ve specified the inputs and ratio of interest, we can start specifying some desirable constraints on the function F that we’re trying to find. The authors speculate that, since vector spaces are linear structures, we should have F encode the information of the ratio in the vector space via vector differences:

$$F(w_i - w_j, \tilde{w}_k) = \frac{P_{ik}}{P_{jk}} \quad (7)$$

which basically says “our representation of the word vectors should be s.t. the **relative** probability of some word \tilde{w}_k occurring in the context of a word w_i compared to it occurring in the context of a different word w_j can be captured by $w_i - w_j$ and \tilde{w}_k alone.”

3. Next we notice that F maps arguments in \mathbb{R}^d to a scalar in \mathbb{R} . The most straightforward way of doing so while maintaining the linear structure we are trying to capture is via a dot product:

$$F((w_i - w_j)^T \tilde{w}_k) = \frac{P_{ik}}{P_{jk}} \quad (8)$$

Note that now $F : \mathbb{R} \mapsto \mathbb{R}$.

4. We want our model to be invariant under the exchanges $w \leftrightarrow \tilde{w}$ and $X \leftrightarrow X^T$. We can restore this symmetry by first requiring that F be a **homomorphism**³ between the groups $(\mathbb{R}, +)$ and $(\mathbb{R}_{>0}, \times)$ (in our case, negation and division would be better symbols, but it’s equivalent). This requires that the following relation hold⁴

$$F(w_i^T \tilde{w}_k - w_j^T \tilde{w}_k) = \frac{F(w_i^T \tilde{w}_k)}{F(w_j^T \tilde{w}_k)} \quad (9)$$

where I’ve grouped terms on the LHS to emphasize how this is the definition of homomorphism. The solution for this equation is that $F(\cdot) \equiv \exp(\cdot)$. Combining this realization with equations 8 and 9 yields,

$$F(w_i^T \tilde{w}_k) = e^{w_i^T \tilde{w}_k} = P_{ik} \quad (10)$$

$$\Rightarrow w_i^T \tilde{w}_k = \log(P_{ik}) = \log(X_{ik}) - \log(X_i) \quad (11)$$

²Pennington et al., “GloVe: Global Vectors for Word Representation.”

³In more detail, $F : (\mathbb{R}, +) \mapsto (\mathbb{R}_{>0}, \times)$, which reads “the function F maps elements in \mathbb{R} and any summation of elements in \mathbb{R} to elements in \mathbb{R} that are greater than zero or any product of positive elements in \mathbb{R} .”

⁴Note that we do not need to know anything about the RHS of the equations above to state this relation. We write it by definition of homomorphism.

where X_{ik} is the number of times word k appears in the context of word i , and $X_i = \sum_k X_{ik}$ is the number of times any word appears in the context of i .

5. Restore symmetry under the exchanges $w \leftrightarrow \tilde{w}$ and $X \leftrightarrow X^T$. We absorb X_i into a bias b_i for w_i since it is independent of k .

$$w_i^T \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik}) \quad (12)$$

6. A main drawback to this model is that it weighs all co-occurrences equally, even those that happen rarely or never. The authors propose a new weighted least squares regression model, introducing a weighting function $f(X_{ij})$ into the cost function of our model:

$$J = \sum_{i,j=1}^V f(X_{ij})(w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij})^2 \quad (13)$$

SPEECH AND LANGUAGE PROCESSING

CONTENTS

2.1	Introduction (Ch. 1 2nd Ed.)	10
2.2	Morphology (Ch. 3 2nd Ed.)	11
2.3	N-Grams (Ch. 6 2nd Ed.)	12
2.4	Naive Bayes and Sentiment (Ch. 6 3rd Ed.)	14
2.5	Hidden Markov Models (Ch. 9 3rd Ed.)	16
2.6	POS Tagging (Ch. 10 3rd Ed.)	19
2.7	Formal Grammars (Ch. 11 3rd Ed.)	22
2.8	Vector Semantics (Ch. 15)	23
2.9	Semantics with Dense Vectors (Ch. 16)	26
2.10	Information Extraction (Ch. 21 3rd Ed)	29

Introduction (Ch. 1 2nd Ed.)

Overview. Going to rapidly jot down what seems most important from this chapter.

- **Morphology:** captures information about the shape and behavior of words in context (Ch. 2/3).
- **Syntax:** knowledge needed to order and group words together.
- **Lexical semantics:** knowledge of the meanings of individual words.
- **Compositional semantics:** knowledge of how these components (words) combine to form larger meanings.
- **Pragmatics:** the appropriate use of polite and indirect language.
- The knowledge of language needed to engage in complex language behavior can be separated into the following 6 distinct categories:
 1. Phonetics and Phonology – The study of linguistic sounds.
 2. Morphology – The study of the meaningful components of words.
 3. Syntax – The study of the structural relationships between words.
 4. Semantics – The study of meaning.
 5. Pragmatics – The study of how language is used to accomplish goals.
 6. Discourse – The study of linguistic units larger than a single utterance.
- Methods for **resolving ambiguity:** pos-tagging, word sense disambiguation, probabilistic parsing, and speech act interpretation.
- **Models and Algorithms.** Among the most important are **state space search** and **dynamic programming** algorithms.

Morphology (Ch. 3 2nd Ed.)

English Morphology. Morphology is the study of the way words are built up from smaller meaning-bearing units, **morphemes**. A morpheme is often defined as the minimal meaning-bearing unit in a language⁵. The two classes of morphemes are **stems** (the “main” morpheme of the word) and **affixes** (the “additional” meanings of various kinds).

Affixes are further divided into prefixes (precede stem), suffixes (follow stem), circumfixes (do both), and infixes (inside the stem).

Two classes of ways to form words from morphemes: **inflection** and **derivation**. Inflection is the combination of a word stem with a grammatical morpheme, usually resulting in a word of the same class as the original stem, and usually filling some syntactic function like agreement. Derivation is the combination of a word stem with a grammatical morpheme, usually resulting in a word of a different class, often with a meaning hard to predict exactly.

⁵Examples: “fox” is its own morpheme, while “cats” consists of the morpheme “cat” and the morpheme “-s”.

Counting Words. Most N -gram based systems deal with the *wordform*, meaning they treat words like “cats” and “cat” as distinct. However, we may want to treat the two as instances of a single abstract word, or **lemma**: a set of lexical forms having the same stem, the same major part of speech, and the same word-sense.

Simple (Unsmoothed) N-Grams. An N -gram is a $N-1$ th order Markov model (because it looks $N-1$ steps in the past). Notation: the authors use the convention that $w_1^n \triangleq w_1, w_2, \dots, w_n$ to denote a sequence of n words. Given this, we can write the general equation for the N -gram approximation for the probability of the n th word ($n > N$) in a sentence:

$$\Pr [w_n | w_1^{n-1}] \approx \Pr [w_n | w_{n-N+1}^{n-1}] \quad (14)$$

for $N \geq 2$. We can compute these probabilities by simply counting:

$$\Pr [w_n | w_1^{n-1}] = \frac{C(w_{n-N+1}^{n-1} w_n)}{C(w_1^{n-1})} \quad (15)$$

where $C(\cdot)$ is the number of times the sequence, denoted as \cdot , occurred in the corpus.

Entropy. Denote the random variable of interest as \mathbf{x} with probability function $p(\mathbf{x})$. The entropy of this random variable is:

$$H(\mathbf{x}) = - \sum_x p(\mathbf{x} = x) \log_2 p(\mathbf{x} = x) \quad (16)$$

which should be thought of as a lower bound on the number of bits it would take to encode a certain decision or piece of information in the optimal coding scheme. The value 2^H is the **perplexity**, which can be interpreted as the weighted average number of choices a random variable has to make.

Cross Entropy for Comparing Models. Useful when we don't know the actual probability distribution p that generated some data. Assume we have some model m that's an approximation of p . The cross-entropy of m on p is defined by:

$$H(p, m) = \lim_{n \rightarrow \infty} \frac{1}{n} \sum_{W \in L} p(w_1, \dots, w_n) \log m(w_1, \dots, w_n) \quad (17)$$

That is we draw sequences according to the probability distribution p , but sum the log of their probability according to m .

Naive Bayes and Sentiment (Ch. 6 3rd Ed.)

Table of Contents Local

Written by Brandon McKinzie

[\[3rd Ed.\]](#) [\[Quick Review\]](#)

Overview. This chapter is concerned with **text categorization**, the task of classifying an entire text by assigning it a label drawn from some set of labels. *Generative classifiers* like naive Bayes build a model of each class. Given an observation, they return the class most likely to have generated the observation. *Discriminative classifiers* like logistic regression instead learn what features from the input are most useful to discriminate between the different possible classes. Notation: we will assume we have a training set of N documents, each hand-labeled with some class: $\{(d_1, c_1), \dots, (d_N, c_N)\}$.

Discriminative systems are often more accurate and hence more commonly used.

Naive Bayes. A multinomial⁶ classifier with a naive assumption about how the features interact. We model a text document as a **bag of words**, meaning we store (1) the words that occurred and (2) their frequencies. It's a probabilistic classifier, meaning it estimates the label/class of a document d as

$$\hat{c} = \arg \max_{c \in C} \Pr [c | d] \quad (19)$$

$$= \arg \max_{c \in C} \frac{\Pr [d | c] \Pr [c]}{\Pr [d]} \quad (20)$$

$$= \arg \max_{c \in C} \Pr [d | c] \Pr [c] \quad (21)$$

Computing the class-conditional distribution (the likelihood) $\Pr [d | c]$ over all possible $d \in D$ is typically intractable; we must introduce some simplifying assumptions and use an approximation of it. Our assumptions in this section will be:

- **Bag-of-Words:** Assume that word position within a document is irrelevant. (Counts still matter)

6

Rapid review: **multinomial distribution.** Let $\mathbf{x} = (x_1, \dots, x_k)$ denote the result of an experiment with n independent trials ($n = \sum_i x_i$) and k possible outcomes for any given trial. i.e. x_i is the number of trials that had outcome i ($1 \leq i \leq k$). The pmf of this multinomial distribution, over all possible \mathbf{x} constrained by $n = \sum_i x_i$, is:

$$\Pr [\mathbf{x} = (x_1, \dots, x_k); n] = \frac{n!}{x_1! \cdots x_k!} p_1^{x_1} \times \cdots \times p_k^{x_k} \quad (18)$$

where p_i is the probability of outcome i for any single trial.

- **Naive Bayes Assumption:** First, recall that d is typically modeled as a (random) vector consisting of features f_1, \dots, f_n , each of which has an associated probability distribution $\Pr [f_i | c]$. The NB assumption is that the features are mutually independent given the class c :

$$\Pr [f_1, \dots, f_n | c] = \Pr [f_1 | c] \cdots \Pr [f_n | c] \quad (22)$$

$$c_{NB} = \arg \max_{c \in C} \Pr [c] \prod_{f \in F} \Pr [f | c] \quad (23)$$

where 23 is the final equation for the class chosen by the naive Bayes classifier.

In text classification we typically use the word at position i in the document as f_i , and move to log space to avoid underflow/increase speed:

$$c_{NB} = \arg \max_{c \in C} \log \Pr [c] + \sum_i^{\text{len}(d)} \log \Pr [w_i | c] \quad (24)$$

Classifiers that use a linear combination of the inputs to make a classification decision (e.g. NB, logistic regression) are called linear classifiers.

Training the NB Classifier. No real "training" in my opinion, just simple counting from the data:

$$\hat{P}[c] = \frac{N_c}{N_{docs}} \quad (25)$$

$$\hat{P}[w_i | c] = \frac{\text{count}(w_i, c) + 1}{(\sum_{w \in V} \text{count}(w, c)) + |V|} \quad (26)$$

The Laplace smoothing is added to avoid the occurrence of zero-probabilities in equation 23.

Optimizing for Sentiment Analysis.

- It often improves performance [for sentiment] to clip word counts in each document to 1 ("binary NB").
- deal with *negations* in some way.
- Use sentiment lexicons, lists of words that are pre-annotated with positive or negative sentiment.

Hidden Markov Models (Ch. 9 3rd Ed.)

Overview. Here we will first go over the math behind HMMs: the **Viterbi**, **Forward**, and **Baum-Welch** (EM) algorithms for unsupervised or semi-supervised learning. Recall that a HMM is defined by specifying the set of N states Q , transition matrix A , sequence of T observations O , sequence of observation likelihoods B , and the initial/final states. They can be characterized by three fundamental problems:

1. **Likelihood.** Given an HMM $\lambda = (A, B)$ and observation sequence, compute the likelihood (prob. of the observations given the model). (**Forward**)
2. **Decoding.** Given an HMM $\lambda = (A, B)$ and observation sequence, discover the best hidden state sequence. (**Viterbi**)
3. **Learning.** Given an observation sequence and the set of states in the HMM, learn the HMM parameters A and B . (**Baum-Welch/Forward-Backward/EM**)

The Forward Algorithm. For likelihood computation. We want to compute the probability of some sequence of observations O , without knowing the sequence of hidden states (that emitted the observations) Q . In general, this can be expressed by summing over all possible hidden state sequences:

$$\Pr [O] = \sum_Q \Pr [Q] \Pr [O | Q] \quad (27)$$

However, for N hidden states and T observations, this summation involves N^T terms, which becomes intractable rather quickly. Instead, we can use the $\mathcal{O}(N^2T)$ **Forward Algorithm**. The forward algorithm can be defined by initialization, recursion definition, and termination, shown respectively as follows:

$$\alpha_1(j) = a_{0j}b_j(o_1) \quad 1 \leq j \leq N \quad (28)$$

$$\alpha_t(j) = \sum_{i=1}^N \alpha_{t-1}(i)a_{ij}b_j(o_t) \quad 1 \leq j \leq N, 1 \leq t \leq T \quad (29)$$

$$\Pr [O | \lambda] = \alpha_T(q_F) = \sum_{i=1}^N \alpha_T(i)a_{iF} \quad (30)$$

Viterbi Algorithm. For decoding. Want the most likely hidden state sequence given observations. Let $v_t(j)$ represent the probability that we are in state j after t observations and passing through the most probable state sequence q_0, q_1, \dots, q_{t-1} . Similar to the forward algorithm, we show the defining equations for the Viterbi algorithm below:

$$v_1(j) = a_{0j}b_j(o_1) \quad 1 \leq j \leq N \quad (31)$$

$$v_t(j) = \max_{i=1}^N v_{t-1}(i)a_{ij}b_j(o_t) \quad (32)$$

$$P^* = v_T(q_F) = \max_{i=1}^N v_T(i)a_{iF} \quad (33)$$

N.B.: At each step, the best path up to that point can be found by taking the *argmax* instead of max.

Baum-Welch Algorithm. AKA forward-backward algorithm, a special case of the EM algorithm. Given an observation sequence O and the set of best possible states in the HMM, learn the HMM parameters A and B . First, we must define some new notation. The **backward probability** β is defined as:

$$\beta_t(i) \triangleq \Pr [o_{t+1}, \dots, o_T \mid q_t = i, \lambda] \quad (34) \quad \text{Remember } \lambda \equiv (A, B)$$

As usual, we can compute its values inductively as follows:

$$\beta_T(i) = a_{iF} \quad 1 \leq i \leq N \quad (35)$$

$$\beta_t(i) = \sum_{j=1}^N a_{ij}b_j(o_{t+1})\beta_{t+1}(j) \quad 1 \leq i \leq N, 1 \leq t \leq T \quad (36)$$

$$\Pr [O \mid \lambda] = \alpha_T(q_F) = \beta_1(q_0) \quad (37)$$

$$= \sum_{j=1}^N a_{0j}b_j(o_1)\beta_1(j) \quad (38)$$

We can use the forward and backward probabilities α and β to compute the transition probabilities a_{ij} and observation probabilities $b_i(o_t)$ from an observation sequence. The derivation steps are as follows:

1. **Estimating \hat{a}_{ij} .** Begin by defining quantities that will prove useful:

$$\xi_t(i, j) \triangleq \Pr [q_t = i, q_{t+1} = j \mid O, \lambda] \quad (39)$$

$$\tilde{\xi}_t(i, j) \triangleq \Pr [q_t = i, q_{t+1} = j, O \mid \lambda] \quad (40)$$

$$= \alpha_t(i) a_{ij} b_j(t+1) \beta_{t+1}(j) \quad (41)$$

Remember, knowing the observation sequence does NOT give us the sequence of hidden states.

where you should be able to derive eq. 41 in your head using just logic. If you cannot, review before continuing. We can then derive $\xi_t(i, j)$ using basic definitions of conditional probability, combined with eq. 37. Finally, we estimate \hat{a}_{ij} as the expected number of transitions $q_i \rightarrow q_j$ divided by the expected number of transitions from q_i total:

$$\hat{a}_{ij} = \frac{\sum_{t=1}^{T-1} \xi_t(i, j)}{\sum_{t=1}^{T-1} \sum_{k=1}^N \xi_t(i, k)} \quad (42)$$

2. **Estimating $\hat{b}_j(v_k)$.** We define our estimate as the expected number of times we are in q_j and emit observation v_k , divided by the expected number of times we are in state j . Similar to our approach for \hat{a}_{ij} we define helper quantities for these values at a given timestep, then sum over them (all t) to obtain our estimate.

$$\gamma_t(j) \triangleq \Pr [q_t = j \mid O, \lambda] \quad (43)$$

$$= \frac{\Pr [q_t = j, O \mid \lambda]}{\Pr [O \mid \lambda]} \quad (44)$$

$$= \frac{\alpha_t(j) \beta_t(j)}{\Pr [O \mid \lambda]} \quad (45)$$

Thus, we obtain $\hat{b}_j(v_k)$ by summing over all timesteps where $o_t = v_k$, denoted as the set T_{v_k} , divided by the summation over all t regardless of o_t :

$$\hat{b}_j(v_k) = \frac{\sum_{t \in T_{v_k}} \gamma_t(j)}{\sum_{t=1}^T \gamma_t(j)} \quad (46)$$

At last we can finally define the **Forward-Backward Algorithm** as follows:

1. Initialize A and B .
2. **E-step.** Compute $\gamma_t(j)$ ($\forall t, j$), and compute $\xi_t(i, j)$ ($\forall t, i, j$).
3. **M-step.** Update all \hat{a}_{ij} and $\hat{b}_j(v_k)$.
4. Upon convergence, return A and B .

POS Tagging (Ch. 10 3rd Ed.)

English Parts-of-Speech. POS are traditionally defined based on syntactic and morphological function, grouping words that have similar neighboring words or take similar affixes.

	Part of Speech	Definition	Properties
Open classes:	noun	people, places, things	occur with determiners, take possessives,
	verb	actions, processes	3rd-person-sg, progressive, past participle
	adjective	properties, qualities	
	adverb	modify verbs, adverbs, verb phrases	directional, locative, degree, manner, tempo

Common nouns can be divided into *count* (e.g. goat/goats) and *mass* (e.g. snow) nouns.

Closed classes. POS with relatively fixed membership. Some of the most important in English are:

prepositions: on, under, over, near, by, at, from, to, with
determiners: a, an, the
pronouns: she, who, I, others
conjunctions: and, but, or, as, if, when
auxiliary verbs: can, may, should, are
particles: up, down, on, off, in, out, at, by
numerals: one, two, three, first, second, third

Some subtleties: the **particle** resembles a preposition or an adverb and is used in combination with a verb. An example case where “over” is a particle: “she turned the paper over.” When a verb and a particle behave as a single syntactic and/or semantic unit, we call the combination a **phrasal verb**. Phrasal verbs cause widespread problems with NLP because they often behave as a semantic unit with a noncompositional meaning – one that is not predictable from the distinct meanings of the verb and the particle. Thus, “turn down” means something like “reject”, “rule out” means “eliminate”, “find out” is “discover”, and “go on” is “continue”.

HMM POS Tagging. Since we typically train on labeled data, we need only use the Viterbi algorithm for decoding⁷. In the POS case, we wish to find the sequence of n tags, \hat{t}_1^n , given the observation sequence of n words w_1^n .

$$\hat{t}_1^n = \arg \max_{t_1^n} \Pr [t_1^n | w_1^n] = \arg \max_{t_1^n} \Pr [w_1^n | t_1^n] \Pr [t_1^n] \quad (47)$$

where we've dropped the denominator after using Bayes' rule (since argmax is the same). HMM taggers made two further simplifying assumptions:

$$\Pr [w_1^n | t_1^n] \approx \prod_{i=1}^n \Pr [w_i | t_i] \quad (48)$$

$$\Pr [t_1^n] \approx \prod_{i=1}^n \Pr [t_i | t_{i-1}] \quad (49)$$

We can thus plug-in these values into eq. 47 to obtain the equation for \hat{t}_1^n .

$$\hat{t}_1^n = \arg \max_{t_1^n} \prod_{i=1}^n \Pr [w_i | t_i] \Pr [t_i | t_{i-1}] \quad (50)$$

$$= \arg \max_{t_1^n} \prod_{i=1}^n b_i(w_i) a_{i-1,i} \quad (51)$$

where I've written a "translated" version on the second line using the familiar syntax from the previous chapter. In practice, we can obtain quick estimates for the two probabilities on the RHS by taking counts/averages over our tagged training data. We then run through the Viterbi algorithm to find all the argmaxes over states for the most likely hidden state sequence.

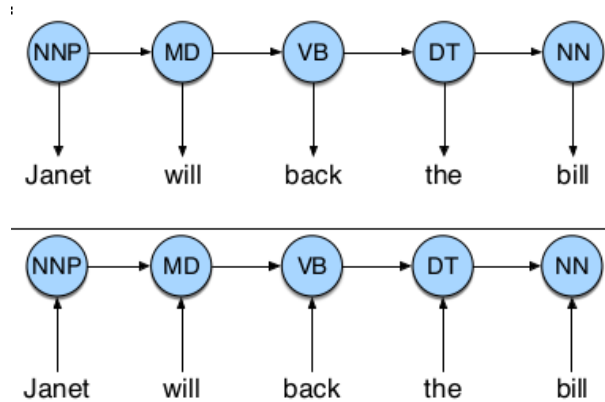
Maximum Entropy Markov Models (MEMMs). A sequence model adaptation of the MaxEnt (multinomial logistic regression) classifier⁸. Since HMMs are generative models, they decompose $\Pr [T | W]$ into $\Pr [W | T] \Pr [T]$ when computing the best tag sequence \hat{T} . Since MEMMs are discriminative, they compute/model $\Pr [T | W]$ directly:

$$\hat{T} = \arg \max_T \Pr [T | W] = \arg \max_T \prod_i \Pr [t_i | w_i, t_{i-1}] \quad (52)$$

Visually, we can think of the difference between HMMs and MEMMs via the direction of arrows, as illustrated below.

⁷Recall that decoding is the problem of finding the best hidden state sequence, given $\lambda = (A, B)$ and observation sequence O .

⁸Because it is based on logistic regression, the MEMM is a **discriminative sequence model**. By contrast, the HMM is a **generative sequence model**.



The top shows the HMM representation, while the bottom is MEMM.

The reason to use a discriminative sequence model is that discriminative models make it easier to incorporate a much wider variety of features.

Bidirectionality. The one problem with the MEMM and HMM models as presented is that they are exclusively run left-to-right. MEMMs⁹ have a weakness known as the **label bias problem**. Consider the tagged fragment: “*will/NN to/TO fight/VB*”¹⁰. Even though the word “*will*” is followed by “*to*”, which strongly suggests “*will*” is a NN, a MEMM will incorrectly label “*will*” as MD (modal verb). The culprit lies in the fact that $\Pr [TO | to, t_{will}]$ is essentially 1 regardless of t_{will} ; i.e. the fact that “*to*” must have the tag TO has **explained away** the presence of TO and so the model doesn’t learn the importance of the previous NN tag for predicting TO.

One way to implement bidirectionality (and thus allowing e.g. the link between TO being available when tagging the NN) is to use a **Conditional Random Field** (CRF) model. However, CRFs are much more computationally expensive than MEMMs and don’t work better for tagging.

⁹And other non-generative finite-state models based on next-state classifiers

¹⁰Note on the tag meanings: TO literally means “to”. MD means “modal” and refers to modal verbs such as *will, shall*, etc.

Formal Grammars (Ch. 11 3rd Ed.)

Constituency and CFGs. Discovering the inventory of constituents present in the language. Groups of words like *noun phrases* or *prepositional phrases* can be thought of as single units which can only be placed within certain parts of a sentence.

The most widely used formal system for modeling constituent structure in English is the **Context-Free Grammar**¹¹. A CFG consists of a set of **productions** (rules), e.g.

$$\text{NP} \longrightarrow \text{Det Nominal} \quad (53)$$

$$\text{NP} \longrightarrow \text{ProperNoun} \quad (54)$$

$$\text{Nominal} \longrightarrow \text{Noun} \mid \text{Nominal Noun} \quad (55)$$

where the arrow is to be read “is composed of” or “consists of.”

The sequence of rule expansions going from left to right is called a **derivation** of the string of words, commonly represented by a **parse tree**. The formal language defined by a CFG is the set of strings that are derivable from the designated **start symbol**.

¹¹Also called Phrase-Structure Grammars. Equiv formalism as Backus-Naur Form (BNF)

Vector Semantics (Ch. 15)

Words and Vectors. Vector models are generally based on a **co-occurrence**, an example of which is a **term-document matrix**: Each row is identified by a word, and each column a document. A given cell value is the number of times the assoc. word occurred in the assoc. document. Can also view each column as a document vector.

Information Retrieval: task of finding document d , from D docs total, that best matches a query q .

For individual word vectors, however, it is most common to instead use a **term-term** matrix¹², in which columns are also identified by individual words. Now, cell values are the number of times the row (target) word and the column (context) word co-occur in some context in some training corpus. The context is most commonly a window around the row/target word, meaning a cell gives the number of times the column word occurs in a window of $\pm N$ words from the row word.

- **Q:** What about the co-occurrence of a word with itself (row i , col i)?
 - **A:** It is included, yes. Source: “The size of the window ... is generally between 1 and 8 words on each side of the target word (*for a total context of 3-17 words*).”
- **Q:** Why is the size of each vector generally $|V|$ (vocab size)? Shouldn't this vary substantially with window and corpus size?
 - **A:** idk

(**TODO:** revisit end of page 5 in my pdf of this)

Pointwise Mutual Information (PMI). Motivation: raw frequencies in a co-occurrence matrix aren't that great, and words like “the” (which aren't useful and occur everywhere) can really skew things. *The best weighting or measure of association between words should tell us how much more often than chance the two words co-occur.* PMI is such a measure.

$$\text{[Mutual Information]} \quad I(X, Y) = \sum_x \sum_y P(X = x, Y = y) \text{PMI}(x, y) \quad (56)$$

$$\text{[PMI]} \quad \text{PMI}(x, y) = \ln \frac{P(x, y)}{P(x)P(y)} \quad (57)$$

which can be applied for our specific use case as $\text{PMI}(w, c) = \ln \frac{P(w, c)}{P(w)P(c)}$. The interpretation is simple: the denominator tells the joint probability of the given target word w occurring with context word c if they were independent of each other, while the numerator tells us how often we observed the two words together (assuming we compute probability by using the MLE).

¹²Also called the word-word or term-context matrix

Therefore, the ratio gives us how an estimate of how much more the target and feature co-occur than we expect by chance¹³. Most people use **Positive PMI**, which is just $\max(0, \text{PMI})$. We can compute a **PPMI matrix** (to replace our co-occurrence matrix), where PPMI_{ij} gives the PPMI value of word w_i with context c_j . The authors show a few formulas which is really distracting since all we actually need is the counts $f_{ij} = \text{counts}(w_i, c_j)$, and from there we can use basic probability and Bayes rule to get the PPMI formula.

- **Q:** Explain why the following is true: very rare words tend to have very high PMI values.
 - **A:** hi
- **Q:** What is the range of α used in PPMI_α ? What is the intuition behind doing this?
 - **A:** For reference:

$$\text{PPMI}_\alpha(w, c) = \max\left(\ln \frac{P(w, c)}{P(w)P_\alpha(c)}, 0\right) \quad (58)$$

$$P_\alpha(c) = \frac{\text{count}(c)^\alpha}{\sum_{c'} \text{count}(c')^\alpha} \quad (59)$$

Although there are better methods than PPMI for weighted co-occurrence matrices, most notably **TF-IDF**, things like tf-idf are not generally used for measuring *word similarity*. For that, PPMI and significance-testing metrics like t-test and likelihood-ratio are more common. The **t-test** statistic, like PMI, measures how much more frequent the association is than chance.

$$t = \frac{\bar{x} - \mu}{\sqrt{s^2/N}} \quad (60)$$

$$\text{t-test}(a, b) = \frac{P(a, b) - P(a)P(b)}{\sqrt{P(a)P(b)}} \quad (61)$$

where \bar{x} is the observed mean, while μ is the expected mean [under our null-hypothesis of independence].

Measuring Similarity. By far the most common similarity metric is the **cosine** of the angle between the vectors:

$$\text{cosine}(\mathbf{v}, \mathbf{w}) = \frac{\mathbf{v} \cdot \mathbf{w}}{|\mathbf{v}| |\mathbf{w}|} \quad (62)$$

Note that, since we've been defining vector elements as frequencies/PPMI values, they won't have negative elements, and thus our cosine similarities will be between 0 and 1 (not -1 and 1).

¹³Computing PMI this way can be problematic for word pairs with small probability, especially if we have a small corpus. Recognize that PMI should never really be negative, but in practice this happens for such cases

Alternatives to cosine:

- **Jaccard measure:** Described as "weighted number of overlapping features, normalized", but looks like a silly hack in my opinion:

$$\text{sim}_{Jac}(\mathbf{v}, \mathbf{w}) = \frac{\sum_{i=1}^N \min(\mathbf{v}_i, \mathbf{w}_i)}{\sum_{i=1}^N \max(\mathbf{v}_i, \mathbf{w}_i)} \quad (63)$$

- **Dice measure:** Another hack. This displeases me.

$$\frac{2 \times \sum_{i=1}^N \min(\mathbf{v}_i, \mathbf{w}_i)}{\sum_{i=1}^N (\mathbf{v}_i + \mathbf{w}_i)} \quad (64)$$

- **Jensen-Shannon Divergence:** An alternative to the KL-divergence¹⁴, which represents the divergence of each distribution from the mean of the two:

$$\text{sim}_{JS}(\mathbf{v}||\mathbf{w}) = D\left(\mathbf{v} \left\| \frac{\mathbf{v} + \mathbf{w}}{2}\right.\right) + D\left(\mathbf{w} \left\| \frac{\mathbf{v} + \mathbf{w}}{2}\right.\right) \quad (66)$$

¹⁴ Idea: if two vectors, \mathbf{v} and \mathbf{w} , each express a probability distribution (their values sum to one), then they are as similar to the extent that these probability distributions are similar. The basis of comparing two probability distributions P and Q is the **Kullback-Leibler** divergence or relative entropy, defined as:

$$D(P||Q) = \sum_x P(x) \log \frac{P(x)}{Q(x)} \quad (65)$$

Semantics with Dense Vectors (Ch. 16)

Table of Contents Local

Written by Brandon McKinzie

Overview. This chapter introduces three methods for generating short, dense vectors: (1) dimensionality reduction like SVD, (2) neural networks like skip-gram or CBOW, and (3) Brown clustering.

Dense Vectors via SVD. Method for finding more important dimensions of a dataset, “important” defined as dimensions wherein the data most varies. First applied (for language) for generating embeddings from term-document matrices in a model called **Latent Semantic Analysis** (LSA). LSA is just SVD on a $|V| \times c$ term-document matrix \mathbf{X} , factorized into $\mathbf{W}\mathbf{\Sigma}\mathbf{C}^T$. By using only the top $k < m$ dimensions of these three matrices, the product becomes a least-squares approx. to the original \mathbf{X} . It also gives us the reduced $|V| \times k$ matrix \mathbf{W}_k , where each row (word) is a k -dimensional vector (embedding). Voilà, we have our dense vectors!

$$\begin{aligned}\mathbf{W} &\in \mathbb{R}^{|V| \times m} \\ \mathbf{\Sigma} &\in \mathbb{R}^{m \times m} \\ \mathbf{C}^T &\in \mathbb{R}^{m \times c}\end{aligned}$$

Note that LSA implementations typically use a particular weighting of each cell in the term-document matrix called the **local** and **global** weights.

$$[\text{local}] \quad \log f(i, j) + 1 \tag{67}$$

$$[\text{global}] \quad 1 + \frac{\sum_j p(i, j) \log p(i, j)}{\log D} \tag{68}$$

$f(i, j)$ is the raw frequency of word i in context j . D is number of docs.

For the case of a word-word matrix, it is common to use PPMI weighting.

Skip-Gram and CBOW. Neural models learn an embedding by starting with a random vector and then iteratively shifting a word’s embeddings to be more like the embeddings of neighboring words, and less like the embeddings of words that don’t occur nearby¹⁵ Word2vec, for example, learns embeddings by training to predict neighboring words¹⁶.

- **Skip-Gram:** Learns two embeddings for each word w : the **word embedding** v (within matrix \mathbf{W}) and **context embedding** c (within matrix \mathbf{C}). Visually:

$$\mathbf{W} = \begin{pmatrix} \mathbf{v}_0^T \\ \mathbf{v}_1^T \\ \vdots \\ \mathbf{v}_{|V|}^T \end{pmatrix} \quad \mathbf{C} = \begin{pmatrix} \mathbf{c}_0 & \mathbf{c}_1 & \cdots & \mathbf{c}_{|V|} \end{pmatrix} \tag{69}$$

¹⁵Why? Why is this a sensible assumption? I see no reason a priori why it ought to be true.

¹⁶Note that the prediction task is not the goal – it just happens to result in good word embeddings. Hacky.

For a context window of $L = 2$, and at a given word $\mathbf{v}^{(t)}$ inside the corpus¹⁷, our goal is to predict the context [words] denoted as $[\mathbf{c}^{(t-2)}, \mathbf{c}^{(t-1)}, \mathbf{c}^{(t+1)}, \mathbf{c}^{(t+2)}]$.

- Example: Consider one of the context words, say $\mathbf{c}^{(t+1)} \triangleq \mathbf{c}_k$, where we also assume it's the k th word in our vocab. Also assume that our target word $\mathbf{v}^{(t)} \triangleq \mathbf{v}_j$ is the j th word in our vocab.
- Our task is to compute $\Pr[\mathbf{c}_k | \mathbf{v}_j]$. We do this with a softmax:

$$\Pr[\mathbf{c}_k | \mathbf{v}_j] = \frac{e^{\mathbf{c}_k^T \mathbf{v}_j}}{\sum_{i \in |V|} e^{\mathbf{c}_i^T \mathbf{v}_j}} \quad (70)$$

- **CBOW**: Continuous bag of words. Basically the mirror-image of skip-gram. Goal is to predict current word $\mathbf{v}^{(t)}$ from the context window of $2L$ words $[\mathbf{c}^{(t-2)}, \mathbf{c}^{(t-1)}, \mathbf{c}^{(t+1)}, \mathbf{c}^{(t+2)}]$.

As usual, the denominator of the softmax is computationally expensive, and usually we approximate it with **negative sampling**.

In the training phase, the algorithm walks through the corpus, at each target word choosing the surrounding context words as positive examples, and for each positive example also choosing k noise samples or negative samples: non-neighbor words. The goal will be to move the embeddings toward the neighbor words and away from the noise words.

Example: Suppose we come along the following window (in “[]”) ($L=2$) in our corpus:

lemon, a [tablespoon of apricot preserves or] jam

Ultimately, we want dot products, $\mathbf{c}_i \cdot \text{vector}(\text{“apricot”})$, to be high for all four of the context words \mathbf{c}_i . We do negative sampling by sampling k random noise words according to their [unigram] frequency. So here, for e.g. $k = 2$, this would amount to 8 noise words, 2 for each context word. We want the dot products between “apricot” and these noise words to be low. For a given single context-word pair (w, c) , our training objective is to maximize:

$$\log \sigma(c \cdot w) + \sum_{i=1}^k \mathbb{E}_{w_i \sim p(w)} [\log \sigma(-w_i \cdot w)] \quad (71)$$

In practice, common to use $p^{3/4}(w)$ instead of $p(w)$

Again, the above is for a single context-target word-pair and, accordingly, the summation is only over $k = 2$ (for our example). Don't try to split the expectation into a summation or anything – just view it as an expected value. To iteratively shift parameters, we use an optimizer like SGD.

¹⁷Note that, technically, the position t of $\mathbf{v}^{(t)}$ is irrelevant for our computation; we are predicting those words based on which word $\mathbf{v}^{(t)}$ is in the vocabulary, not it's position in the corpus.

The actual model architecture is a typical neural net, progressing as follows:

$$\text{“apricot”} \rightarrow \mathbf{w}^{\text{one-hot}} = [0 \ 0 \ \dots \ 1 \ \dots \ 0] \quad (72)$$

$$\rightarrow \mathbf{h} = \mathbf{W}^T \mathbf{w}^{\text{one-hot}} \quad (73)$$

$$\rightarrow \mathbf{o} = \mathbf{C}^T \mathbf{h} = [c_0^T \mathbf{h}, c_1^T \mathbf{h}, \dots, c_{|V|}^T \mathbf{h}]^T \quad (74)$$

$$\rightarrow \mathbf{y} = \text{softmax}(\mathbf{o}) = [\text{Pr}[c_0|\mathbf{h}], \text{Pr}[c_2|\mathbf{h}], \dots, \text{Pr}[c_{|V|}|\mathbf{h}]]^T \quad (75)$$

$$(76)$$

Brown Clustering. An agglomerative clustering algorithm for deriving vector representations of words by clustering words based on their associations with the preceding or following words. Makes use of the **class-based language model** (CBLM), wherein each word w belongs to some class $c \in C$ via the probability $P(w | c)$. CBLMs define

$$P(w_i | w_{i-1}) = P(c_i | c_{i-1})P(w_i | c_i) \quad (77)$$

$$P(\text{corpus} | C) = \prod_{i=1}^n P(w_i | w_{i-1}) \quad (78)$$

A naive and extremely inefficient version of Brown clustering, a hierarchical clustering algorithm, is as follows:

1. Each word is initially assigned to its own cluster.
2. For each cluster pair $(c_i, c_{j \neq i})$, compute the value of eq 78 that would result from merging c_i and c_j into a single cluster. The pair whose merger results in the smallest decrease in eq 78 is merged.
3. Clustering proceeds until all words are in one big cluster.

This process builds a binary tree from the bottom-up, and the binary string corresp. to a word’s traversal from leaf-to-root is its representation.

Information Extraction (Ch. 21 3rd Ed)

Overview. The first step in most IE tasks is **named entity recognition** (NER). Next we can do **relation extraction**: finding and classifying semantic relations among the entities, e.g. “spouse-of.” **Event extraction** is finding the events in which the entities participate, and **event coreference** for figuring out which event mentions actually refer to the same event. It’s also common to extract dates/times (temporal expression) and perform **temporal expression normalization** to map them onto specific calendar dates. Finally, we can do **template filling**: finding recurring/stereotypical situations in documents and filling the template slots with appropriate material.

Named Entity Recognition. Standard algorithm is word-by-word sequence labeling task by a MEMM or CRF, trained to label tokens with tags indicating the presence of particular kinds of NEs. It is common to label with **BIO tagging**, for beginning, inside, and outside of entities. If we have n unique entity types, then we’d have $2n + 1$ BIO tags¹⁸. A helpful illustration is shown below:

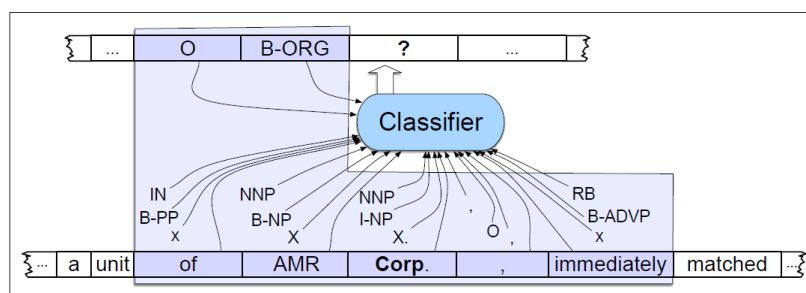


Figure 21.7 Named entity recognition as sequence labeling. The features available to the classifier during training and classification are those in the boxed area.

Here we see a classifier determining the label for *Corp.* with a context window of size 2 and various features shown in the boxed region. For evaluation of NER, we typically use the familiar **recall**, **precision**, and **F1 measure**.

¹⁸ $2n$ for B-<NE> and I-<NE>, with +1 for the blanket *O* tag (not any of our NEs)

Relation Extraction. The four main algorithm classes used are (1) hand-written patterns, (2) supervised ML, (3) semi-supervised, and (4) unsupervised. Terminology:

- **Infobox:** structured tables associated with certain articles/topics/etc. For example, the Wikipedia infobox for Stanford includes structured facts like `state = 'California'`.
- **Resource Description Framework (RDF):** a metalanguage of RDF triples, tuples of (entity, relation, entity), called a subject-predicate-object expression. For example: (Golden Gate Park, location, San Francisco).
- **hyponym:** the “is-a” relation.
- **hyponym:** the “kind-of” relation. *Gelidium is a kind of red algae.*

Overview of the four algorithm classes:

1. **Patterns.** Consider a sentence that has the following form:

$$NP_0 \text{ such as } NP_1\{, NP_2, \dots, (and|or)NP_i\}, i \geq 1$$

also known as a *lexico-syntactic pattern*, which implies $\forall NP_i, i \geq 1, \text{hyponym}(NP_i, NP_0)$ ¹⁹. Patterns typically have high precision but low-recall.

2. **Supervised.** The general approach for finding relations in a given sequence of words is the following:

- (a) Find all pairs of named entities in the sequence (typically a single sentence).
- (b) For each pair, use a trained binary classifier to predict whether or not the entities in the pair are indeed related.
- (c) If related, use a classifier trained to predict the relation given the entity-pair.

As with NER, **the most important step** in this process is to identify useful surface features that will be useful for relation classification, including word features, NER features, syntactic paths (chunk seqs, constituent paths, dependency-tree paths), and more.

3. **Semi-supervised** via bootstrapping. Suppose we have a few high-precision **seed patterns** (or seed tuples)²⁰. **Bootstrapping** proceeds by taking the entities in the seed pair, and then finding sentences (on the web, or whatever dataset we are using) that contain both entities. From all such sentences, we extract and generalize the context around the entities to learn new patterns.

```

function BOOTSTRAP(Relation R) returns new relation tuples

    tuples ← Gather a set of seed tuples that have relation R
    iterate
        sentences ← find sentences that contain entities in seeds
        patterns ← generalize the context between and around entities in sentences
        newpairs ← use patterns to grep for more tuples
        newpairs ← newpairs with high confidence
        tuples ← tuples + newpairs
    return tuples

```

Figure 21.14 Bootstrapping from seed entity pairs to learn relations.

¹⁹Here, $\text{hyponym}(A, B)$ means “A is a kind-of (hyponym) of B.”

²⁰seed tuples are tuples of the general form (M1, M2) where M1 and M2 are each specific named entities we know have the relation of interest R.

4. **Unsupervised.** The Re Verb system extracts a relation from a sentence s in 4 steps:

1. Run a part-of-speech tagger and entity chunker over s
2. For each verb in s , find the longest sequence of words w that start with a verb and satisfy syntactic and lexical constraints, merging adjacent matches.
3. For each phrase w , find the nearest noun phrase x to the left which is not a relative pronoun, wh-word or existential “there”. Find the nearest noun phrase y to the right.
4. Assign confidence c to the relation $r = (x, w, y)$ using a confidence classifier and return it.

Event Extraction. An event mention is any expression denoting an event or state that can be assigned to a particular point, or interval, in time. Note that this is quite different than the colloquial usage of the word “event,” you should think of the two as distinct. Here, most event mentions correspond to verbs, and most verbs introduce events. Event extraction is typically modeled via ML, detecting events via sequence models with BIO tagging, and assigning event classes/attributes with multi-class classifiers.

Template Filling. The task is creation of one template for each event in the input documents, with the slots filled with text from the document. For example, an event could be “Fare-Raise Attempt” with corresponding template (slots to be filled) “(<Lead Airline>, <Amount>, <Effective Date>, <Follower>)”. This is generally modeled by training two separate supervised systems:

1. **Template recognition.** Trained to determine if template T is present in sentence S . Here, “present” means there is a sequence within the sentence that could be used to fill a slot within template T .
2. **Role-filler extraction.** Trained to detect each role (slot-name), e.g. “Lead Airline”.