

# 4 N-GRAMS

*But it must be recognized that the notion “probability of a sentence” is an entirely useless one, under any known interpretation of this term.*

Noam Chomsky (1969, p. 57)

*Anytime a linguist leaves the group the recognition rate goes up.*

Fred Jelinek (then of the IBM speech group) (1988)<sup>1</sup>

Radar O’Reilly, the mild-mannered clerk of the 4077th M\*A\*S\*H unit, had an uncanny ability to guess the next word someone was going to say. In this chapter we take up this idea of word prediction; what word, for example, is likely to follow:

I’d like to make a collect...

Hopefully most of you concluded that a very likely word is *call*, or *international* or *phone*, but probably not *the*. We formalize this idea of **word prediction** with probabilistic models called ***N*-grams**, which predict the next word from the previous  $N - 1$  words. Such statistical models of word sequences are also called **language models** or **LMs**. Computing the probability of the next word will turn out to be closely related to computing the probability of a sequence of words. The following sequence, for example, has a non-zero probability of appearing in a text:

... all of a sudden I notice three guys standing on the sidewalk...

while this same set of words in a different order has a very low probability:

on guys all I of notice sidewalk three a sudden standing the

As we will see, estimators like *N*-grams that assign a conditional probability to possible next words can be used to assign a joint probability to an entire sentence.

<sup>1</sup> This wording from his address is as recalled by Jelinek himself; the quote didn’t appear in the proceedings (Palmer and Finin, 1990). Some remember a more snappy version: *Every time I fire a linguist the performance of the recognizer improves.*

WORD PREDICTION

LANGUAGE MODELS

LM

Whether estimating probabilities of next words or of whole sequences, the  $N$ -gram model is one of the most important tools in speech and language processing.

$N$ -grams are essential in any task in which we have to identify words in noisy, ambiguous input. In **speech recognition**, for example, the input speech sounds are very confusable and many words sound extremely similar. Russell and Norvig (1995) give an intuition from **handwriting recognition** for how probabilities of word sequences can help. In the movie *Take the Money and Run*, Woody Allen tries to rob a bank with a sloppily written hold-up note that the teller incorrectly reads as “I have a gub”. Any speech and language processing system could avoid making this mistake by using the knowledge that the sequence “I have a gun” is far more probable than the non-word “I have a gub” or even “I have a gull”.

$N$ -gram models are also essential in statistical **machine translation**. Suppose we are translating a Chinese source sentence 他向记者介绍了该声明的主要内容 and as part of the process we have a set of potential rough English translations:

he briefed to reporters on the chief contents of the statement  
 he briefed reporters on the chief contents of the statement  
 he briefed to reporters on the main contents of the statement  
**he briefed reporters on the main contents of the statement**

An  $N$ -gram grammar might tell us that, even after controlling for length, *briefed reporters* is more likely than *briefed to reporters*, and *main contents* is more likely than **chief contents**. This lets us select the bold-faced sentence above as the most fluent translation sentence, i.e. the one that has the highest probability.

In **spelling correction**, we need to find and correct spelling errors like the following (from Kukich (1992)) that accidentally result in real English words:

They are leaving in about fifteen *minuets* to go to her house.  
 The design *an* construction of the system will take more than a year.

Since these errors have real words, we can't find them by just flagging words not in the dictionary. But note that *in about fifteen minuets* is a much less probable sequence than *in about fifteen minutes*. A spellchecker can use a probability estimator both to detect these errors and to suggest higher-probability correction.

Word prediction is also important for **augmentative communication** (Newell et al., 1998) systems that help the disabled. People who are unable to use speech or sign-language to communicate, like the physicist Steven Hawking, can communicate by using simple body movements to select words from a menu that are spoken by the system. Word prediction can be used to suggest likely words for the menu.

Besides these sample areas,  $N$ -grams are also crucial in NLP foundations like **part-of-speech tagging**, **natural language generation**, and **word similarity**, as well as applications from **authorship identification** and **sentiment extraction** to **predictive text input** systems for cell phones.

## 4.1 COUNTING WORDS IN CORPORA

[upon being asked if there weren't enough words in the English language for him]:

*"Yes, there are enough, but they aren't the right ones."*

James Joyce, reported in Bates (1997)

CORPUS  
CORPORA

Probabilities are based on counting things. Before we talk about probabilities, we need to decide what we are going to count. Counting of things in natural language is based on a **corpus** (plural **corpora**), an on-line collection of text or speech. Let's look at two popular corpora, Brown and Switchboard. The Brown Corpus is a 1 million word collection of samples from 500 written texts from different genres (newspaper, novels, non-fiction, academic, etc.), assembled at Brown University in 1963-64 (Kučera and Francis, 1967; Francis, 1979; Francis and Kučera, 1982). How many words are in the following Brown sentence?

(4.1) He stepped out into the hall, was delighted to encounter a water brother.

Example (4.1) has 13 words if we don't count punctuation-marks as words, 15 if we count punctuation. Whether we treat period ("."), comma (","), and so on as words depends on the task. Punctuation is critical for finding boundaries of things (comma, periods, colons), and for identifying some aspects of meaning (question marks, exclamation marks, quotation marks). For some tasks, such as part-of-speech tagging or parsing or sometimes speech synthesis, we thus sometimes treat punctuation as if they were separate words.

UTTERANCE

The Switchboard corpus of telephone conversations between strangers was collected in the early 1990s and contains 2430 conversations averaging 6 minutes each, totaling 240 hours of speech and about 3 million words (Godfrey et al., 1992). Such corpora of spoken language don't have punctuation, but do introduce other complications with defining words. Let's look at one utterance from Switchboard; an **utterance** is the spoken correlate of a sentence:

(4.2) I do uh main- mainly business data processing

DISFLUENCIES  
FRAGMENT  
FILLERS  
FILLED PAUSES

This utterance has two kinds of **disfluencies**. The broken-off word *main-* is called a **fragment**. Words like *uh* and *um* are called **fillers** or **filled pauses**. Should we consider these to be words? Again, it depends on the application. If we are building an automatic dictation system based on automatic speech recognition, we might want to eventually strip out the disfluencies.

But we also sometimes keep disfluencies around. How disfluent a person is can be used to identify them, or to detect if they are stressed or confused. Disfluencies also often occur with particular syntactic structures, so they may help in parsing and word prediction. Stolcke and Shriberg (1996) found for example that treating *uh* as a word improves next-word prediction (why might this be?), and so most speech recognition systems treat *uh* and *um* as words.<sup>2</sup>

Are capitalized tokens like *They* and uncapitalized tokens like *they* the same word? These are lumped together in speech recognition, while for part-of-speech-

<sup>2</sup> Clark and Fox Tree (2002) showed that *uh* and *um* have different meanings. What do you think they are?

tagging capitalization is retained as a separate feature. For the rest of this chapter we will assume our models are not case-sensitive.

How about inflected forms like *cats* versus *cat*? These two words have the same **lemma** *cat* but are different wordforms. Recall from Ch. 3 that a lemma is a set of lexical forms having the same stem, the same major part-of-speech, and the same word-sense. The **wordform** is the full inflected or derived form of the word. For morphologically complex languages like Arabic we often need to deal with lemmatization. *N*-grams for speech recognition in English, however, and all the examples in this chapter, are based on wordforms.

As we can see, *N*-gram models, and counting words in general, requires that we do the kind of tokenization or text normalization that we introduced in the previous chapter; separating out punctuation, dealing with abbreviations like *m.p.h.*, normalizing spelling, and so on.

How many words are there in English? To answer this question we need to distinguish **types**, the number of distinct words in a corpus or vocabulary size  $V$ , from **tokens**, the total number  $N$  of running words. The following Brown sentence has 16 tokens and 14 types (not counting punctuation):

(4.3) They picnicked by the pool, then lay back on the grass and looked at the stars.

The Switchboard corpus has about 20,000 wordform types (from about 3 million wordform tokens) Shakespeare's complete works have 29,066 wordform types (from 884,647 wordform tokens) (Kučera, 1992) The Brown corpus has 61,805 wordform types from 37,851 lemma types (from 1 million wordform tokens). Looking at a very large corpus of 583 million wordform tokens, Brown et al. (1992a) found that it included 293,181 different wordform types. Dictionaries can help in giving lemma counts; dictionary entries, or **boldface forms** are a very rough upper bound on the number of lemmas (since some lemmas have multiple boldface forms). The American Heritage third edition dictionary lists 200,000 boldface forms. It seems like the larger corpora we look at, the more word types we find. In general (?) suggest that the vocabulary size (the number of types) grows with at least the square root of the number of tokens (i.e.  $V > O(\sqrt{N})$ ).

In the rest of this chapter we will continue to distinguish between types and tokens, using "types" to mean wordform types.

## 4.2 SIMPLE (UNSMOOTHED) *N*-GRAMS

Let's start with some intuitive motivations for *N*-grams; we assume that the reader has acquired some very basic background in probability theory. Our goal is to compute the probability of a word  $w$  given some history  $h$ , or  $P(w|h)$ . Suppose the history  $h$  is "*the water is so transparent that*" and we want to know the probability that the next word is *the*:

$$(4.4) \quad P(\text{the}|\text{the water is so transparent that}).$$

How can we compute this probability? One way is to estimate it from relative frequency counts. For example, we could take a very large corpus, count the number of times we

see *the water is so transparent that*, and count the number of times this is followed by *the*. This would be answering the question “Out of the times we saw the history  $h$ , how many times was it followed by the word  $w$ ”, as follows:

$$(4.5) \quad P(\text{the}|\text{the water is so transparent that}) = \frac{C(\text{the water is so transparent that the})}{C(\text{the water is so transparent that})}$$

With a large enough corpus, such as the web, we can compute these counts, and estimate the probability from Equation (4.5). You should pause now, go to the web and compute this estimate for yourself.

While this method of estimating probabilities directly from counts works fine in many cases, it turns out that even the web isn’t big enough to give us good estimates in most cases. This is because language is creative; new sentences are created all the time, and we won’t always be able to count entire sentences. Even simple extensions of the example sentence may have counts of zero on the web (such as “*the water of Walden Pond is so transparent that the*”).

Similarly, if we wanted to know the joint probability of an entire sequence of words like *the water is so transparent*, we could do it by asking “out of all possible sequences of 5 words, how many of them are *the water is so transparent*?” We would have to get the count of *the water is so transparent*, and divide by the sum of the counts of all possible 5 word sequences. That seems rather a lot of work to estimate!

For this reason, we’ll need to introduce cleverer ways of estimating the probability of a word  $w$  given a history  $h$ , or the probability of an entire word sequence  $W$ . Let’s start with a little formalizing of notation. In order to represent the probability of a particular random variable  $X_i$  taking on the value “the”, or  $P(X_i = \text{“the”})$ , we will use the simplification  $P(\text{the})$ . We’ll represent a sequence of  $N$  words either as  $w_1 \dots w_n$  or  $w_1^n$ . For the joint probability of each word in a sequence having a particular value  $P(X = w_1, Y = w_2, Z = w_3, \dots)$  we’ll use  $P(w_1, w_2, \dots, w_n)$ .

Now how can we compute probabilities of entire sequences like  $P(w_1, w_2, \dots, w_n)$ ? One thing we can do is to decompose this probability using the **chain rule of probability**:

$$(4.6) \quad \begin{aligned} P(X_1 \dots X_n) &= P(X_1)P(X_2|X_1)P(X_3|X_1^2) \dots P(X_n|X_1^{n-1}) \\ &= \prod_{k=1}^n P(X_k|X_1^{k-1}) \end{aligned}$$

Applying the chain rule to words, we get:

$$(4.7) \quad \begin{aligned} P(w_1^n) &= P(w_1)P(w_2|w_1)P(w_3|w_1^2) \dots P(w_n|w_1^{n-1}) \\ &= \prod_{k=1}^n P(w_k|w_1^{k-1}) \end{aligned}$$

The chain rule shows the link between computing the joint probability of a sequence and computing the conditional probability of a word given previous words. Equation (4.7) suggests that we could estimate the joint probability of an entire sequence of words by multiplying together a number of conditional probabilities. But using the chain rule doesn't really seem to help us! We don't know any way to compute the exact probability of a word given a long sequence of preceding words,  $P(w_n|w_1^{n-1})$ . As we said above, we can't just estimate by counting the number of times every word occurs following every long string, because language is creative and any particular context might have never occurred before!

The intuition of the  $N$ -gram model is that instead of computing the probability of a word given its entire history, we will **approximate** the history by just the last few words.

**BIGRAM** The **bigram** model, for example, approximates the probability of a word given all the previous words  $P(w_n|w_1^{n-1})$  by the conditional probability of the preceding word  $P(w_n|w_{n-1})$ . In other words, instead of computing the probability

$$(4.8) \quad P(\text{the}|\text{The water of Walden Pond is so transparent that})$$

we approximate it with the probability

$$(4.9) \quad P(\text{the}|\text{that})$$

When we use a bigram model to predict the conditional probability of the next word we are thus making the following approximation:

$$(4.10) \quad P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-N+1}^{n-1})$$

**MARKOV** This assumption that the probability of a word depends only on the previous word is called a **Markov** assumption. Markov models are the class of probabilistic models that assume that we can predict the probability of some future unit without looking too far into the past. We can generalize the bigram (which looks one word into the past) to the trigram (which looks two words into the past) and thus to the **N-gram** (which looks  $N - 1$  words into the past).

Thus the general equation for this  $N$ -gram approximation to the conditional probability of the next word in a sequence is:

$$(4.11) \quad P(w_n|w_1^{n-1}) \approx P(w_n|w_{n-N+1}^{n-1})$$

Given the bigram assumption for the probability of an individual word, we can compute the probability of a complete word sequence by substituting Equation (4.10) into Equation (4.7):

$$(4.12) \quad P(w_1^n) \approx \prod_{k=1}^n P(w_k|w_{k-1})$$

How do we estimate these bigram or  $N$ -gram probabilities? The simplest and most intuitive way to estimate probabilities is called **Maximum Likelihood Estimation**, or **MLE**. We get the MLE estimate for the parameters of an  $N$ -gram models by

taking counts from a corpus, and **normalizing** them so they lie between 0 and 1.<sup>3</sup>

NORMALIZING

For example, to compute a particular bigram probability of a word  $y$  given a previous word  $x$ , we'll compute the count of the bigram  $C(xy)$  and normalize by the sum of all the bigrams that share the same first word  $x$ :

$$(4.13) \quad P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{\sum_w C(w_{n-1}w)}$$

We can simplify this equation, since the sum of all bigram counts that start with a given word  $w_{n-1}$  must be equal to the unigram count for that word  $w_{n-1}$ . (The reader should take a moment to be convinced of this):

$$(4.14) \quad P(w_n|w_{n-1}) = \frac{C(w_{n-1}w_n)}{C(w_{n-1})}$$

Let's work through an example using a mini-corpus of three sentences. We'll first need to augment each sentence with a special symbol  $\langle s \rangle$  at the beginning of the sentence, to give us the bigram context of the first word. We'll also need a special end-symbol  $\langle /s \rangle$ .<sup>4</sup>

```

<s> I am Sam </s>
<s> Sam I am </s>
<s> I do not like green eggs and ham </s>

```

Here are the calculations for some of the bigram probabilities from this corpus

$$\begin{aligned}
 P(I | \langle s \rangle) &= \frac{2}{3} = .66 & P(\text{Sam} | \langle s \rangle) &= \frac{1}{3} = .33 & P(\text{am} | I) &= \frac{2}{3} = .33 \\
 P(\langle /s \rangle | \text{Sam}) &= \frac{1}{2} = 0.5 & P(\langle s \rangle | \text{Sam}) &= \frac{1}{2} = 0.5 & P(\text{Sam} | \text{am}) &= \frac{1}{2} = .5 \\
 P(\text{do} | I) &= \frac{1}{3} = .33
 \end{aligned}$$

For the general case of MLE  $N$ -gram parameter estimation:

$$(4.15) \quad P(w_n|w_{n-N+1}^{n-1}) = \frac{C(w_{n-N+1}^{n-1}w_n)}{C(w_{n-N+1}^{n-1})}$$

RELATIVE  
FREQUENCY

Equation 4.15 (like equation 4.14) estimates the  $N$ -gram probability by dividing the observed frequency of a particular sequence by the observed frequency of a prefix. This ratio is called a **relative frequency**; the use of relative frequencies as a way to estimate probabilities is one example of MLE. In Maximum Likelihood Estimation, the resulting parameter set maximizes the likelihood of the training set  $T$  given the model  $M$  (i.e.,  $P(T|M)$ ). For example, suppose the word *Chinese* occurs 400 times in a corpus of a million words like the Brown corpus. What is the probability that a random word selected from some other text of say a million words will be the word *Chinese*? The MLE estimate of its probability is  $\frac{400}{1000000}$  or .0004. Now .0004 is not

<sup>3</sup> For probabilistic models, normalizing means dividing by some total count so that the resulting probabilities fall legally between 0 and 1.

<sup>4</sup> As Chen and Goodman (1998) point out, we need the end-symbol to make the bigram grammar a true probability distribution. Without an end-symbol, the sentence probabilities for all sentences of a given length would sum to one, and the probability of the whole language would be infinite.

the best possible estimate of the probability of *Chinese* occurring in all situations; it might turn out that in some OTHER corpus or context *Chinese* is a very unlikely word. But it is the probability that makes it *most likely* that Chinese will occur 400 times in a million-word corpus. We will see ways to modify the MLE estimates slightly to get better probability estimates in Sec. 4.5.

Let's move on to some examples from a slightly larger corpus than our 14-word example above. We'll use data from the now-defunct Berkeley Restaurant Project, a dialogue system from the last century that answered questions from a database of restaurants in Berkeley, California (Jurafsky et al., 1994). Here are some sample user queries, lowercased and with no punctuation; a representative corpus of 9332 sentences is on the website:

```
can you tell me about any good cantonese restaurants close by
mid priced thai food is what i'm looking for
tell me about chez panisse
can you give me a listing of the kinds of food that are available
i'm looking for a good place to eat breakfast
when is caffe venezia open during the day
```

Fig. 4.1 shows the bigram counts from a piece of a bigram grammar from the Berkeley Restaurant Project. Note that the majority of the values are zero. In fact, we have chosen the sample words to cohere with each other; a matrix selected from a random set of seven words would be even more sparse.

	i	want	to	eat	chinese	food	lunch	spend
i	5	827	0	9	0	0	0	2
want	2	0	608	1	6	6	5	1
to	2	0	4	686	2	0	6	211
eat	0	0	2	0	16	2	42	0
chinese	1	0	0	0	0	82	1	0
food	15	0	15	0	1	4	0	0
lunch	2	0	0	0	0	1	0	0
spend	1	0	1	0	0	0	0	0

**Figure 4.1** Bigram counts for eight of the words (out of  $V = 1446$ ) in the Berkeley Restaurant Project corpus of 9332 sentences.

Fig. 4.2 shows the bigram probabilities after normalization (dividing each row by the following appropriate unigram counts):

i	want	to	eat	chinese	food	lunch	spend
2533	927	2417	746	158	1093	341	278

Here are a few other useful probabilities:

$$P(i | <s>) = 0.25 \quad P(\text{english} | \text{want}) = 0.0011$$

$$P(\text{food} | \text{english}) = 0.5 \quad P(</s> | \text{food}) = 0.68$$

	i	want	to	eat	chinese	food	lunch	spend
i	0.002	0.33	0	0.0036	0	0	0	0.00079
want	0.0022	0	0.66	0.0011	0.0065	0.0065	0.0054	0.0011
to	0.00083	0	0.0017	0.28	0.00083	0	0.0025	0.087
eat	0	0	0.0027	0	0.021	0.0027	0.056	0
chinese	0.0063	0	0	0	0	0.52	0.0063	0
food	0.014	0	0.014	0	0.00092	0.0037	0	0
lunch	0.0059	0	0	0	0	0.0029	0	0
spend	0.0036	0	0.0036	0	0	0	0	0

**Figure 4.2** Bigram probabilities for eight words in the Berkeley Restaurant Project corpus of 9332 sentences.

Now we can compute the probability of sentences like *I want English food* or *I want Chinese food* by simply multiplying the appropriate bigram probabilities together, as follows:

$$\begin{aligned}
 &P(\langle s \rangle \text{ i want english food } \langle /s \rangle) \\
 &= P(\text{i} | \langle s \rangle)P(\text{want} | \text{i})P(\text{english} | \text{want}) \\
 &\quad P(\text{food} | \text{english})P(\langle /s \rangle | \text{food}) \\
 &= .25 \times .33 \times .0011 \times 0.5 \times 0.68 \\
 &= .000031
 \end{aligned}$$

We leave it as an exercise for the reader to compute the probability of *i want chinese food*. But that exercise does suggest that we'll want to think a bit about what kinds of linguistic phenomena are captured in bigrams. Some of the bigram probabilities above encode some facts that we think of as strictly syntactic in nature, like the fact that what comes after *eat* is usually a noun or an adjective, or that what comes after *to* is usually a verb. Others might be more cultural than linguistic, like the low probability of anyone asking for advice on finding English food.

Although we will generally show bigram models in this chapter for pedagogical purposes, note that when there is sufficient training data we are more likely to use **trigram** models, which condition on the previous two words rather than the previous word. To compute trigram probabilities at the very beginning of sentence, we can use two pseudo-words for the first trigram (i.e.,  $P(\text{I} | \langle s \rangle \langle s \rangle)$ ).

### 4.3 TRAINING AND TEST SETS

The  $N$ -gram model is a good example of the kind of statistical models that we will be seeing throughout speech and language processing. The probabilities of an  $N$ -gram come from the corpus it is trained on. In general, the parameters of a statistical model are trained on some set of data, and then we apply the models to some new data in some task (such as speech recognition) and see how well they work. Of course this new data or task won't be the exact same data we trained on.

TRAINING SET

TEST SET

We can formalize this idea of training on some data, and testing on some other data by talking about these two data sets as a **training set** and a **test set** (or a **training corpus** and a **test corpus**). Thus when using a statistical model of language given some corpus of relevant data, we start by dividing the data into training and test sets. We train the statistical parameters of the model on the training set, and then use this trained model to compute probabilities on the test set.

EVALUATE

This training-and-testing paradigm can also be used to **evaluate** different  $N$ -gram architectures. Suppose we want to compare different language models (such as those based on  $N$ -grams of different order  $N$ , or using the different **smoothing** algorithms to be introduced in Sec. 4.5). We can do this by taking a corpus and dividing it into a training set and a test set. Then we train the two different  $N$ -gram models on the training set and see which one better models the test set. But what does it mean to “model the test set”? There is a useful metric for how well a given statistical model matches a test corpus, called **perplexity**, introduced on page-13. Perplexity is based on computing the probabilities of each sentence in the test set; intuitively, whichever model assigns a higher probability to the test set (hence more accurately predicts the test set) is a better model.

Since our evaluation metric is based on test set probability, it’s important not to let the test sentences into the training set. Suppose we are trying to compute the probability of a particular “test” sentence. If our test sentence is part of the training corpus, we will mistakenly assign it an artificially high probability when it occurs in the test set. We call this situation **training on the test set**. Training on the test set introduces a bias that makes the probabilities all look too high, and causes huge inaccuracies in perplexity.

HELD-OUT

In addition to training and test sets, other divisions of data are often useful. Sometimes we need an extra source of data to augment the training set. Such extra data is called a **held-out** set, because we hold it out from our training set when we train our  $N$ -gram counts. The held-out corpus is then used to set some other parameters; for example we will see the use of held-out data to set interpolation weights in **interpolated**  $N$ -gram models in Sec. 4.6 Finally, sometimes we need to have multiple test sets. This happens because we might use a particular test set so often that we implicitly tune to its characteristics. Then we might need a fresh test set which is truly unseen. In such cases, we call the initial test set the **development** test set or, **devset**. We will discuss development test sets again in Ch. 5.

DEVELOPMENT

How do we divide our data into training, dev test, and test sets? There is a tradeoff, since we want our test set to be as large as possible, since a small test set may be accidentally unrepresentative. On the other hand, we want as much training data as possible. At the minimum, we would want to pick the smallest test set that gives us enough statistical power to measure a statistically significant difference between two potential models. In practice, we often just divide our data into 80% training, 10% development, and 10% test. Given a large corpus that we want to divide into training and test, test data can either be taken from some continuous sequence of text inside the corpus, or we can remove smaller “stripes” of text from randomly selected parts of our corpus and combine them into a test set.

### 4.3.1 N-gram Sensitivity to the Training Corpus

The  $N$ -gram model, like many statistical models, is very dependent on the training corpus. One implication of this is that the probabilities often encode very specific facts about a given training corpus. Another implication is that  $N$ -grams do a better and better job of modeling the training corpus as we increase the value of  $N$ .

We can visualize both of these facts by borrowing the technique of Shannon (1951) and Miller and Selfridge (1950), of generating random sentences from different  $N$ -gram models. It's simplest to visualize how this works for the unigram case. Imagine all the words of English covering the probability space between 0 and 1, each word covering an interval of size equal to its frequency. We choose a random number between 0 and 1, and print out the word whose interval includes the real value we have chosen. We continue choosing random numbers and generating words until we randomly generate the sentence-final token  $\langle /s \rangle$ . The same technique can be used to generate bigrams by first generating a random bigram that starts with  $\langle s \rangle$  (according to its bigram probability), then choosing a random bigram to follow it (again, according to its conditional probability), and so on.

To give an intuition for the increasing power of higher order  $N$ -grams, Fig. 4.3 shows random sentences generated from unigram, bigram, trigram, and quadrigram models trained on Shakespeare's works.

The longer the context on which we train the model, the more coherent the sentences. In the unigram sentences, there is no coherent relation between words, nor sentence-final punctuation. The bigram sentences have some very local word-to-word coherence (especially if we consider that punctuation counts as a word). The trigram and quadrigram sentences are beginning to look a lot like Shakespeare. Indeed a careful investigation of the quadrigram sentences shows that they look a little too much like Shakespeare. The words *It cannot be but so* are directly from *King John*. This is because, not to put the knock on Shakespeare, his oeuvre is not very large as corpora go ( $N = 884,647$ ,  $V = 29,066$ ), and our  $N$ -gram probability matrices are ridiculously sparse. There are  $V^2 = 844,000,000$  possible bigrams alone, and the number of possible quadrigrams is  $V^4 = 7 \times 10^{17}$ . Thus once the generator has chosen the first quadrigram (*It cannot be but*), there are only five possible continuations (*that, I, he, thou, and so*); indeed for many quadrigrams there is only one continuation.

To get an idea of the dependence of a grammar on its training set, let's look at an  $N$ -gram grammar trained on a completely different corpus: the Wall Street Journal (WSJ) newspaper. Shakespeare and the Wall Street Journal are both English, so we might expect some overlap between our  $N$ -grams for the two genres. In order to check whether this is true, Fig. 4.4 shows sentences generated by unigram, bigram, and trigram grammars trained on 40 million words from WSJ.

Compare these examples to the pseudo-Shakespeare in Fig. 4.3. While superficially they both seem to model "English-like sentences" there is obviously no overlap whatsoever in possible sentences, and little if any overlap even in small phrases. This stark difference tells us that statistical models are likely to be pretty useless as predictors if the training sets and the test sets are as different as Shakespeare and WSJ.

Unigram	<ul style="list-style-type: none"> <li>• To him swallowed confess hear both. Which. Of save on trail for are ay device and rote life have</li> <li>• Every enter now severally so, let</li> <li>• Hill he late speaks; or! a more to leg less first you enter</li> <li>• Are where exeunt and sighs have rise excellency took of.. Sleep knave we. near; vile like</li> </ul>
Bigram	<ul style="list-style-type: none"> <li>• What means, sir. I confess she? then all sorts, he is trim, captain.</li> <li>• Why dost stand forth thy canopy, forsooth; he is this palpable hit the King Henry. Live king. Follow.</li> <li>• What we, hath got so she that I rest and sent to scold and nature bankrupt, nor the first gentleman?</li> <li>• Enter Menenius, if it so many good direction found'st thou art a strong upon command of fear not a liberal largess given away, Falstaff! Exeunt</li> </ul>
Trigram	<ul style="list-style-type: none"> <li>• Sweet prince, Falstaff shall die. Harry of Monmouth's grave.</li> <li>• This shall forbid it should be branded, if renown made it empty.</li> <li>• Indeed the duke; and had a very good friend.</li> <li>• Fly, and will rid me these news of price. Therefore the sadness of parting, as they say, 'tis done.</li> </ul>
Quadrigram	<ul style="list-style-type: none"> <li>• King Henry. What! I will go seek the traitor Gloucester. Exeunt some of the watch. A great banquet serv'd in;</li> <li>• Will you not tell me who I am?</li> <li>• It cannot be but so.</li> <li>• Indeed the short and the long. Marry, 'tis a noble Lepidus.</li> </ul>
<p><b>Figure 4.3</b> Sentences randomly generated from four <math>N</math>-grams computed from Shakespeare's works. All characters were mapped to lower-case and punctuation marks were treated as words. Output is hand-corrected for capitalization to improve readability.</p>	

*unigram*: Months the my and issue of year foreign new exchange's september were recession exchange new endorsed a acquire to six executives

*bigram*: Last December through the way to preserve the Hudson corporation N. B. E. C. Taylor would seem to complete the major central planners one point five percent of U. S. E. has already old M. X. corporation of living on information such as more frequently fishing to keep her

*trigram*: They also point to ninety nine point six billion dollars from two hundred four oh six three percent of the rates of interest stores as Mexico and Brazil on market conditions

**Figure 4.4** Sentences randomly generated from three orders of  $N$ -gram computed from 40 million words of the Wall Street Journal. All characters were mapped to lower-case and punctuation marks were treated as words. Output is hand-corrected for capitalization to improve readability.

How should we deal with this problem when we build  $N$ -gram models? In general we need to be sure to use a training corpus that looks like our test corpus. We especially wouldn't choose training and tests from different **genres** of text like newspaper text, early English fiction, telephone conversations, and web pages. Sometimes finding appropriate training text for a specific new task can be difficult; to build  $N$ -grams for text prediction in SMS (Short Message Service), we need a training corpus of SMS data. To build  $N$ -grams on business meetings, we would need to have corpora

of transcribed business meetings.

For general research where we know we want written English but don't have a domain in mind, we can use a balanced training corpus that includes cross-sections from different genres, such as the 1-million word Brown corpus of English (Francis and Kučera, 1982) or the 100-million word British National Corpus (Leech et al., 1994).

Recent research has also studied ways to dynamically **adapt** language models to different genres; see Sec. 4.9.4.

### 4.3.2 Unknown Words: Open versus closed vocabulary tasks

CLOSED  
VOCABULARY

Sometimes we have a language task in which we know all the words that can occur, and hence we know the vocabulary size  $V$  in advance. The **closed vocabulary** assumption is the assumption that we have such a lexicon, and that the test set can only contain words from this lexicon. The closed vocabulary task thus assumes there are no unknown words.

But of course this is a simplification; as we suggested earlier, the number of unseen words grows constantly, so we can't possibly know in advance exactly how many there are, and we'd like our model to do something reasonable with them. We call these unseen events **unknown** words, or **out of vocabulary (OOV)** words. The percentage of OOV words that appear in the test set is called the **OOV rate**.

OOV

OPEN VOCABULARY

An **open vocabulary** system is one where we model these potential unknown words in the test set by adding a pseudo-word called  $\langle \text{UNK} \rangle$ . We can train the probabilities of the unknown word model  $\langle \text{UNK} \rangle$  as follows:

1. **Choose a vocabulary** (word list) which is fixed in advance.
2. **Convert** in the training set any word that is not in this set (any OOV word) to the unknown word token  $\langle \text{UNK} \rangle$  in a text normalization step.
3. **Estimate** the probabilities for  $\langle \text{UNK} \rangle$  from its counts just like any other regular word in the training set.

## 4.4 EVALUATING $N$ -GRAMS: PERPLEXITY

IN VIVO

The correct way to evaluate the performance of a language model is to embed it in an application and measure the total performance of the application. Such end-to-end evaluation, also called **in vivo** evaluation, is the only way to know if a particular improvement in a component is really going to help the task at hand. Thus for speech recognition, we can compare the performance of two language models by running the speech recognizer twice, once with each language model, and seeing which gives the more accurate transcription.

Unfortunately, end-to-end evaluation is often very expensive; evaluating a large speech recognition test set, for example, takes hours or even days. Thus we would like a metric that can be used to quickly evaluate potential improvements in a language model. **Perplexity** is the most common evaluation metric for  $N$ -gram language models. While an improvement in perplexity does not guarantee an improvement in speech recognition performance (or any other end-to-end metric), it often correlates