

1 Chapter 38: Statistical methods for corpus exploitation

2 **Marco Baroni** (baroni@sslmit.unibo.it)
SITLLEC, University of Bologna
Corso Diaz 64, 47100 Forlì, Italy

Stefan Evert (stefan.evert@uos.de)
Institute of Cognitive Science, University of Osnabrück
49069 Osnabrück, Germany

3 Very Final Draft, 20 September 2006

4 **Contents**

5	1 Introduction	1
6	2 The logic behind hypothesis testing	3
7	3 Estimation and effect size	8
8	4 The normal approximation	11
9	5 Two-sample tests	14
10	6 Linguistic units and populations	17
11	7 Non-randomness and the unit of sampling	19
12	8 Other techniques	22
13	9 Directions for further study	24

14 **1 Introduction**

15 Linguists look for generalizations and explanations of various kinds of linguistic phenom-
16 ena. While the interest is usually in an *intensional* view of these phenomena, to be ex-
17 plained in terms of the human language competence, such competence cannot be directly
18 observed. Thus, evidence has to come from an external reflection of it, i.e., it has to be
19 based on an *extensional* view of language. According to this extensional view, a language
20 is defined as the set of all utterances produced by speakers of the language (with all the
21 paradoxes that this view implies – see, e.g., Chomsky 1986, Chapter 2). Corpora are finite
22 samples from the infinite set that constitutes a language in this extensional sense. For
23 example, in this perspective, the Brown corpus (see article 22) is a finite sample of all
24 the utterances produced in written form by American English speakers. Psycholinguis-
25 tic experiments, such as eye-tracking tests, priming, and even traditional grammaticality

1 judgments (Schütze 1996) constitute other sources of evidence. It is important to observe
2 that the empirical analysis of these other sources also requires an extensional view of
3 language.

4 It is rarely the case that linguists are interested in the samples *per se*, rather than
5 in generalizations from the samples to the infinite amount of text corresponding to the
6 extensional definition of a (sub)language. For example, a linguist studying a pattern in the
7 500 text samples of the Brown corpus will typically be interested in drawing conclusions
8 about (written) American English as a whole, and not just about the specific texts that
9 compose the Brown. Statistical inference allows the linguist to generalize from properties
10 observed in a specific sample (corpus) to the same properties in the language as a whole
11 (statistical inference, on the other hand, will not be of help in solving thorny issues such
12 as what is the appropriate extensional definition of a “language as a whole” and how we
13 can sample from that).

14 Statistical inference requires that the problem at hand is *operationalized* in quantita-
15 tive terms, typically in the form of units that can be *counted* in the available sample(s).
16 This is the case we will concentrate on here (but see Section 8 for other kinds of mea-
17 surements). For example, a linguist might be interested in the issue of whether a certain
18 variety of English is more “formal” than another (as in some of Douglas Biber’s work, see
19 article 40). In order to operationalize this research question, the linguist might decide to
20 take passivization as a cue of formality, and count the number of sentences that display
21 passivization in samples from the two varieties. Statistical inference can then be used
22 to generalize from the difference in number of passives between the two samples to the
23 difference between the two varieties that the samples represent (we will discuss this ex-
24 ample and the appropriate techniques further in Section 5). Similarly, a linguist might be
25 interested in whether (certain classes of) idiomatic constructions have a tendency to repel
26 passive formation (as observed by Culicover/Jackendoff 2005 and many others). In order
27 to operationalize this question, the linguist may count the number of passives in idiomatic
28 and non-idiomatic phrases in a corpus. Statistical inference will then help to determine
29 how reliably the attested difference in passive frequency would generalize to idiomatic and
30 non-idiomatic phrases at large. Of course, it is up to the linguist to interpret the gen-
31 eralizations about frequencies produced by statistical analysis in terms of the linguistic
32 phenomena of interest.

33 Statistical inference is necessary because any sample from a language is subject to ran-
34 dom variation. Suppose that someone doubted the claim that non-idiomatic constructions
35 are more prone to passivization than idiomatic constructions, and we wanted to dispel
36 these doubts. A sample of language that reveals a higher proportion of passives among
37 the non-idiomatic constructions, especially if the difference in proportions is small, would
38 not allow us to reject the doubters’ hypothesis: even if they were right, we could not expect
39 the proportions to be exactly identical in *all* samples of language. Statistical inference can
40 help us to determine to what extent the difference between a sample-based observation
41 and a theoretical prediction can be taken as serious evidence that the prediction made
42 by the theory is wrong, and to what extent it can reasonably be attributed to random
43 variation. In the case at hand, statistical inference would tell us whether the difference in
44 passive rates in the two samples can be explained by random variation, or whether it is the
45 symptom of a true underlying difference. It is perhaps worth clarifying from the outset
46 that randomness due to sampling has to be distinguished from measurement errors, such
47 as those introduced by the automatic annotation and analysis of corpus data (something
48 that statistical methods will not help us correct). Suppose that a very skilled linguist sam-
49 pled 100 English sentences and recorded very carefully how many of them are passives,

1 without making any errors. It should be intuitive that, given another random sample of
2 100 sentences and the same error-free linguist, the exact number of passives would proba-
3 bly be different from the one found in the previous sample. This is the random variation
4 we are referring to here.

5 Notice that the necessity of statistical inference pertains to the need to generalize
6 from a finite (random) sample of language data to the theoretically infinite amount of
7 text corresponding to the extensional definition of an entire (sub)language, and it has
8 nothing to do with whether our theory about the phenomenon at hand, or about language
9 competence in general, includes a probabilistic component. The prediction that idiomatic
10 sentences repel the passive construction might stem from a completely categorical theory
11 of how passives and idiomaticity interact – still, randomly sampled English sentences will
12 display a certain amount of variation in the exact proportion of passives they contain.

13 The rest of this article introduces the basics of statistical inference. We use the artifi-
14 cially simple example of testing a hypothesis about the proportion of passives in English
15 sentences (and later proportions of passives in sentences from different English genres), in
16 order to focus on the general philosophy and methodology of statistical inference as ap-
17 plied to corpus linguistics, rather than on the technical details of carrying out the relevant
18 computations, which can be found in many general books on the subject and are imple-
19 mented in all standard statistical packages (see references in Section 9). Section 6 gives
20 examples of how statistical inference can be applied to more realistic linguistic analysis
21 settings.

22 **2 The logic behind hypothesis testing**

23 Imagine that an American English style guide claims that 15% of the sentences in the
24 English language are in the passive voice (as of June 2006, [http://www.ego4u.com/en/
25 business-english/grammar/passive](http://www.ego4u.com/en/business-english/grammar/passive) makes the even bolder statement that no more than
26 10% of English sentences are in the passive voice and writers should be careful to use
27 passives sparingly). This is a fairly easy claim to operationalize, since it is already phrased
28 in terms of a proportion. However, we still need to define what we understand by “the
29 English language”, and what it means for a sentence to be in the passive voice. Given
30 the source of the claim and our need for an extensional definition, it makes sense to
31 take “English” to mean the set of all English texts published in the US and produced by
32 professional writers. Regarding the second issue, we consider a sentence to be in the passive
33 voice if it contains at least one verb in the passive form, which seems to be a plausible
34 interpretation of what the style guide means (after all, it is warning against the overuse of
35 passives), and at the same time makes it easier to count the number of sentences in passive
36 voice using automated pattern matching techniques (which might not be relevant with the
37 small samples we use here, but would be important when dealing with large amounts of
38 data).

39 It is of course impossible to look at all sentences in all the publications satisfying
40 the criteria above – what we can do, at best, is to select a random sample of them. In
41 particular, we took a random sample of 100 sentences of the relevant kind, and we counted
42 the number of them containing a passive. For convenience, we restricted ourselves to
43 publications from 1961, because we are lucky enough to already own a random sample of
44 sentences of the relevant kind from that year – namely, the Brown corpus! All we had
45 to do was select 100 random sentences from this random sample (we will see in Section 7
46 that it is not entirely correct to treat sentences from the Brown as a random sample, but
47 we ignore this for now).

1 If the style guide’s claim is true, we would expect 15 sentences to be in the passive
2 voice. Instead, we found 19 passives. This seems to indicate that the proportion is higher
3 than 15% and rather close to 20%. However, it is obvious that, even if the claim of the
4 style guide was correct, not *all* samples of size 100 would have exactly 15 passives, because
5 of random variation. In light of this, how do we decide whether 19 passives are enough to
6 reject the style guide’s claim?

7 In statistical terms, the claim that we want to verify is called a *null hypothesis*, $H_0 : \pi =$
8 15%, where π is the putative proportion of passives in the set of sentences that constitute
9 our extensional definition of American English. This set of sentences is usually called a
10 *population* in statistical parlance, and the goal of statistical inference is to draw conclusions
11 about certain properties of this population from an available sample (the population itself
12 is practically infinite for all intents and purposes, and we can only access a small finite
13 subset of it). We will often refer to π as a population proportion or parameter in what
14 follows. The number of sentences we have randomly sampled from the population is called
15 the *sample size*, $n = 100$. Intuitively, we expect $e = n \cdot \pi = 15$ passives in the sample if
16 the null hypothesis is true. This is called the *expected frequency*. The number of passives
17 we actually observed in the sample, $f = 19$, is called the *observed frequency*.

18 Having introduced the terminology, we can rephrase the problem above as follows. If
19 we are prepared to reject the null hypothesis that $\pi = 15\%$ for an observation of $f = 19$,
20 there is a certain risk that in doing so we are making the wrong decision. The question is
21 how we can quantify this risk and decide whether it is an acceptable risk to take. Imagine
22 that the null hypothesis in fact holds, and that a large number of linguists perform the
23 same experiment, sampling 100 sentences and counting the passives. We can then formally
24 define risk by the percentage of linguists who wrongly reject the null hypothesis, and thus
25 publish incorrect results. In particular, if our observation of $f = 19$ is deemed sufficient for
26 rejection, all the other linguists who observed 19 or even more passives in their samples
27 would also reject the hypothesis. The risk is thus given by the percentage of samples
28 containing 19 or more passives that would be drawn from a language in which the true
29 proportion of passives is indeed 15%, as stipulated by H_0 . Rejecting the null hypothesis
30 when it is in fact true is known as a *type-1 error* in the technical literature (failure to
31 reject H_0 when it does not hold constitutes a *type-2 error*, which we do not discuss here,
32 but see, e.g., DeGroot/Schervish 2002, Chapter 8).

33 Fortunately, we do not need to hire hundreds of linguists to compute the risk of wrong
34 rejection, since the thought experiment above is fully equivalent to drawing balls from an
35 urn. Each ball represents a sentence of the language, with red balls for passive sentences
36 and white balls for sentences in other voices. The null hypothesis stipulates that the
37 proportion of red balls in the urn is 15%. The observed number of red balls (passives)
38 changes from sample to sample. In statistical terminology, it is called a *random variable*,
39 typically denoted by a capital letter such as X . We simulate a large number of samples
40 from the urn with a computer and tabulate how often each possible value k of the random
41 variable X is observed. The result of this simulation is shown in Figure 1, which reports
42 percentages of samples that yield $X = k$ for k ranging from 0 to 30 (the percentage is
43 indistinguishable from 0 for all values outside this range). For instance, the value $X = 19$
44 can be observed in 5.6% of the samples. The information presented in this graph is called
45 the *sampling distribution* of X under H_0 . The percentage of samples with $X = k$ is called
46 the *probability* $\Pr(X = k)$. For example, $\Pr(X = 19) = 5.6\%$ (our reasoning in this section
47 has led us to what is known as the *frequentist* definition of probabilities; we do not discuss
48 the alternative *Bayesian* interpretation of probability theory here, but see for example
49 Section 1.2 of DeGroot/Schervish 2002).

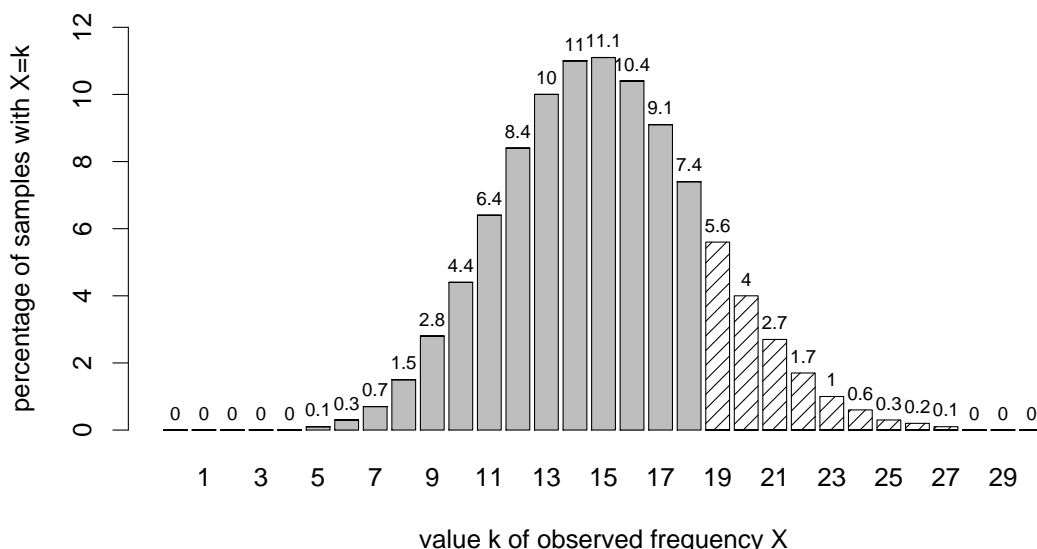


Figure 1: Sampling distribution of X with $n = 100$ and $\pi = 15\%$.

1 Following the discussion above, the risk of wrongly rejecting the null hypothesis for
 2 an observation of $f = 19$ is given by the percentage of samples with $X \geq 19$ in the
 3 sampling distribution, i.e., the probability $\Pr(X \geq 19)$. This probability can be computed
 4 by summing over the shaded bars in Figure 1:

$$\Pr(X \geq 19) = \Pr(X = 19) + \Pr(X = 20) + \dots + \Pr(X = 100) = 16.3\% \quad (1)$$

5 This is called a *tail probability* because it sums over the right-hand “tail” of the dis-
 6 tribution. In the same way, we can compute the risk associated with any other value f ,
 7 namely the probability:

$$\Pr(X \geq f) := \sum_{k=f}^n \Pr(X = k) \quad (2)$$

8 We refer to this risk as the *p-value* of an observation f . Notice that *smaller* p-values
 9 are *more* significant, since they indicate that it is less risky to reject the null hypothesis,
 10 and hence they allow greater confidence in the conclusion that the null hypothesis should
 11 be rejected. In our example, the p-value associated with $f = 19$ indicates that the risk of
 12 false rejection is unacceptably high at $p = 16.3\%$. If we used $f = 19$ as the threshold for
 13 rejection and the null hypothesis happened to be true, about one in six experiments would
 14 lead to wrong conclusions. In order to decide whether to reject H_0 or not, the computed
 15 p-value is often compared with a conventional scale of *significance levels*. A p-value below
 16 5% is always required to consider a result significant. Other common significance levels
 17 are 1% and 0.1% (usually written as mathematical fractions rather than percentages and
 18 denoted by the symbol α , viz. $\alpha = .05$, $\alpha = .01$ and $\alpha = .001$).

19 So far, we have only been considering cases in which the observed frequency is greater
 20 than the one predicted under H_0 – reflecting our intuition that the proportion of passives
 21 proposed by the style guide errs on the side of being too low, rather than too high.
 22 However, in principle it would also be possible that the proportion of passives is *lower*
 23 than predicted by H_0 . Coming back to our passive-counting linguists, if they are prepared

1 to reject H_0 for $f = 19$, they should also reject it for $X = 7$ or $X = 8$, since these values
 2 are even more “extreme” than 19 with respect to H_0 . Thus, when computing the p-value
 3 of $f = 19$, it is typically appropriate to sum over the probabilities of all values that are at
 4 least as “extreme” as the observed value, to either side of the expected frequency e , since
 5 they add to the risk of false rejection.

6 It is difficult to determine exactly which of the values below e should count as equally
 7 extreme as, or more extreme than f , but one reasonable approach is to include all the
 8 values of X with an absolute difference $|X - e| \geq |f - e|$. Using $|f - e|$ (or, more precisely,
 9 $(f - e)^2$, which has certain mathematical advantages) as a measure of “extremeness” leads
 10 to a class of statistical tests known as *chi-squared tests*. An alternative approach would
 11 rather compare the probabilities $\Pr(X = k)$ and $\Pr(X = f)$ as a measure of extremeness,
 12 resulting in a class known as likelihood (ratio) tests. In many common cases, both classes
 13 of tests give very similar results. We will focus on chi-squared tests in this article, but see
 14 article 57 for an application where likelihood tests are known to be superior.

15 In the case at hand, using the chi-squared criterion, the p-value would be computed
 16 by adding up the probabilities of $X \geq 19$ and $X \leq 11$ (since $|19 - 15| = 4 = |11 - 15|$).
 17 In the illustration shown in Figure 1, we would add the bars for $k = 1 \dots 11$ to the shaded
 18 area. This way of computing the p-value, taking both “extreme” tails of the distribution
 19 into account, is called a *two-tailed test* (the approach above, where we considered only one
 20 side, is known as a *one-tailed test*). Of course, a two-tailed p-value is always greater than
 21 (or equal to) the corresponding one-tailed p-value. In our running example, the two-tailed
 22 p-value obtained by summing over the bars for $X \leq 11$ and $X \geq 19$ turns out to be 32.6%,
 23 indicating a very high risk if we chose to reject the null hypothesis for $f = 19$ (you might
 24 obtain a different p-value if the binomial test implemented in your software package uses
 25 the likelihood criterion, although the value will still indicate a very high risk in case of
 26 rejection). Had our experiment yielded 22 passives instead, the one-tailed test would have
 27 produced a p-value of 3.9%, while the two-tailed test would have given a p-value of 6.7%.
 28 Thus, by adopting the common 5% significance threshold, we would have had enough
 29 evidence to reject the null hypothesis according to the one-tailed test, but not enough
 30 according to the two-tailed test.

31 As a general rule, one should always use the more conservative two-tailed test, unless
 32 there are very strong reasons to believe that the null hypothesis could only be violated
 33 in one direction – but it is hard to think of linguistic problems where this is the case
 34 (in many situations we can predict the probable direction of the violation, but there are
 35 very few cases where we would be ready to claim that a violation in the other direction is
 36 absolutely impossible). If we use a two-tailed test, the interpretation of a significant result
 37 will of course have to take into account whether f is greater or smaller than e . Observed
 38 frequencies of 25 and 5 passives, respectively, would both lead to a clear rejection of the
 39 null hypothesis that 15% of all sentences are in the passive voice, but they would require
 40 rather different explanations.

41 Not only have we been spared the expense of hiring passive-counting linguists to repeat
 42 the experiment; it is not even necessary to perform expensive computer simulation exper-
 43 iments in order to carry out the sort of tests we just illustrated, because $\Pr(X = k)$ – the
 44 percentage of samples of size n from a population with proportion π of passive sentences
 45 that would result in a certain value k of the random variable X – can be computed with
 46 the following formula, known as the *binomial distribution* (the hypothesis test we have
 47 described above, unsurprisingly, is called a *binomial test*):

$$\Pr(X = k) = \binom{n}{k} (\pi)^k (1 - \pi)^{n-k} \quad (3)$$

1 The binomial coefficient $\binom{n}{k}$, “ n choose k ”, represents the number of ways in which an
2 unordered set of k elements can be selected from n elements. Any elementary textbook on
3 probability theory or statistics will show how to compute it; see, e.g., DeGroot/Schervish
4 (2002, Section 1.8). Of course, all statistical software packages implement binomial coef-
5 ficients and the binomial distribution.

6 For a different null hypothesis about the population proportion π or a different sample
7 size n , we obtain sampling distributions with different peaks and shapes – in statistical
8 terminology, π and n are the *parameters* of the binomial distribution. In particular, the
9 value of π affects the location of the peak in the histogram. For example, if we hypothesized
10 that $\pi = 30\%$, we would see a peak around the expected value $e = n \cdot \pi = 30$ in the
11 histogram corresponding to Figure 1. Intuitively, experiments in which we draw 1,000
12 balls will tend to produce outcomes that are closer to the expected value than experiments
13 in which we draw 100 balls. Thus, by decreasing or increasing n , we obtain distributions
14 that have narrower or wider shapes, respectively. A sample of size 100 is small by the
15 standards of statistical inference. As Karl Pearson, one of the founding fathers of modern
16 statistics, once put it: “Only naughty brewers deal in small samples!” (cf. Pearson 1990,
17 73; this quip was a reference to W. S. Gosset, an employee of the Guinness brewery who
18 developed and published the now famous t-test under the pseudonym of “Student”). It
19 will typically be difficult to reject H_0 based on such a sample, unless the true proportion
20 is very far away from the null hypothesis, exactly because a small sample size leads to a
21 wide sampling distribution. Had we taken a sample of 1,000 sentences and counted 190
22 passives, the null hypothesis would have been clearly rejected (a two-sided binomial test
23 with $f = 190$, $n = 1000$ and $H_0 : \pi = 15\%$ gives a p-value of $p = 0.048\%$, sufficient for
24 rejection even at the very conservative significance level $\alpha = .001$).

25 The procedure of hypothesis testing that we introduced in this section is fundamental
26 to understanding statistical inference. At the same time, it is not entirely intuitive. Thus,
27 before we move on, we want to summarize its basic steps. For the whole process to be
28 meaningful, we must have a *null hypothesis* H_0 that operationalizes a research question
29 in terms of a quantity that can be computed from observable data. In our case, the null
30 hypothesis stipulates that the proportion of passives in the *population* of (professionally
31 written American) English sentences is 15%, i.e.: $H_0 : \pi = 15\%$. We draw a random
32 *sample* of size n of the relevant units (100 sentences in our case) from the population,
33 and count the number of units that have the property of interest (in our case, being
34 passive sentences). Given the population proportion stipulated by the null hypothesis
35 and the sample size, we can determine a *sampling distribution* (by simulation or using
36 a mathematical formula). The sampling distribution specifies, for each possible outcome
37 of the experiment (expressed by the *random variable* X , which in our case keeps track
38 of the frequency of passives in the sample), how likely it is under the null hypothesis.
39 This *probability* is given by the percentage of a large number of experiments that would
40 produce the outcome X in a world in which the null hypothesis is in fact true. The
41 sampling distribution allows us, for every possible value k of X , to compute the *risk* of
42 making a mistake when we are prepared to reject the null hypothesis for $X = k$. This risk,
43 known as the *p-value* corresponding to k , is given by the overall percentage of experiments
44 that give an outcome at least as extreme as $X = k$ in a world in which the null hypothesis
45 is true (see above for the *one-* and *two-tailed* ways to interpret what counts as “extreme”).
46 At this point, we look at the actual outcome of the experiment in our sample, i.e., the
47 *observed* quantity f (in our case, f is the number of passives in a sample of 100 sentences),
48 and we compute the p-value (risk) associated with f . In our example, the (two-tailed)
49 p-value is 32.6%, indicating a rather high risk in rejecting the null hypothesis. We can

1 compare the p-value we obtained with conventional thresholds, or *significance levels*, that
2 correspond to “socially acceptable” levels of risk, such as the 5% threshold $\alpha = .05$. If
3 the p-value is higher than the threshold, we say that the results of the experiment are not
4 *statistically significant*, i.e., there is a non-negligible possibility that the results would be
5 obtained by chance even if the null hypothesis is true.

6 Notice that a non-significant result simply means that our evidence is not strong enough
7 to reject the null hypothesis. It does *not* tell us that the null hypothesis is correct. In our
8 example, although the observed frequency is not entirely unlikely under the null hypothesis
9 of a passive proportion of 15%, there are many other hypotheses under which the same
10 result would be even more likely, most obviously, the hypothesis that the population pro-
11 portion is 19%. Because of this indirect nature of statistical hypothesis testing, problems
12 undergoing statistical treatment are typically operationalized in a way in which the null
13 hypothesis is “uninteresting”, or contradicts the theory we want to support. Our hope is
14 that the evidence we gather is strong enough to reject H_0 . We will come back to this in
15 Section 5 below, presenting a two-sample setting where this strategy should sound more
16 natural.

17 While many problems require more sophisticated statistical tools than the ones de-
18 scribed in this section, the basic principles of hypothesis testing will be exactly the same
19 as in the example we just discussed.

20 3 Estimation and effect size

21 Suppose that we ran the experiment with a sample of $n = 1,000$ sentences, $f = 190$ of
22 which turned out to be in the passive voice. As we saw in the previous section, this result
23 with the larger sample leads to a clear rejection of the null hypothesis $H_0 : \pi = 15\%$. At
24 this point, we would naturally like to know what the *true* proportion of passives is in edited
25 American English. Intuitively, our best guess is the observed proportion of passives in the
26 sample, i.e., $\hat{\pi} = f/n$. This intuitive choice can also be justified mathematically. It is then
27 known as a *maximum-likelihood estimate* or *MLE* (DeGroot/Schervish 2002, Section 6.5).

28 Since we have estimated a single value for the population proportion, $\hat{\pi}$ is called a *point*
29 *estimate*. The problem with point estimates is that they are subject to the same amount
30 of random variation as the observed frequency on which they are based: most linguists
31 performing the same experiment would obtain a different estimate $\hat{\pi} = X/n$ (note that,
32 mathematically speaking, $\hat{\pi}$ is a random variable just like X , which assumes a different
33 value for each sample).

34 Let us put the question in a slightly different way: besides the point estimate $\hat{\pi} = 19\%$,
35 which other values of π are also plausible given our observation of $f = 190$ passives in a
36 sample of $n = 1,000$ sentences? Since $H_0 : \pi = 15\%$ was rejected by the binomial test, we
37 know for instance that the value $\pi = 15\%$ is *not* plausible according to our observation.
38 This approach allows us to answer the question in an indirect way. For any potential
39 estimate $\pi = x$, we can perform a binomial test with the null hypothesis $H_0 : \pi = x$ in
40 order to determine whether the value x is plausible (H_0 cannot be rejected at the chosen
41 significance level α) or not (H_0 can be rejected). Note that failure to reject H_0 does not
42 imply that the estimate x is very likely to be accurate, but only that we cannot rule out
43 the possibility $\pi = x$ with sufficient confidence. Figure 2 illustrates this procedure for six
44 different values of x , when $f = 190$ and $n = 1,000$. As the figure shows, $H_0 : \pi = 17\%$
45 would not be rejected, and thus 17% is in our set of plausible values. On the other hand,
46 $H_0 : \pi = 16.5\%$ would be rejected, and thus 16.5% is not in our set.

47 Collecting all plausible values $\pi = x$, we obtain a *confidence set*. For the binomial

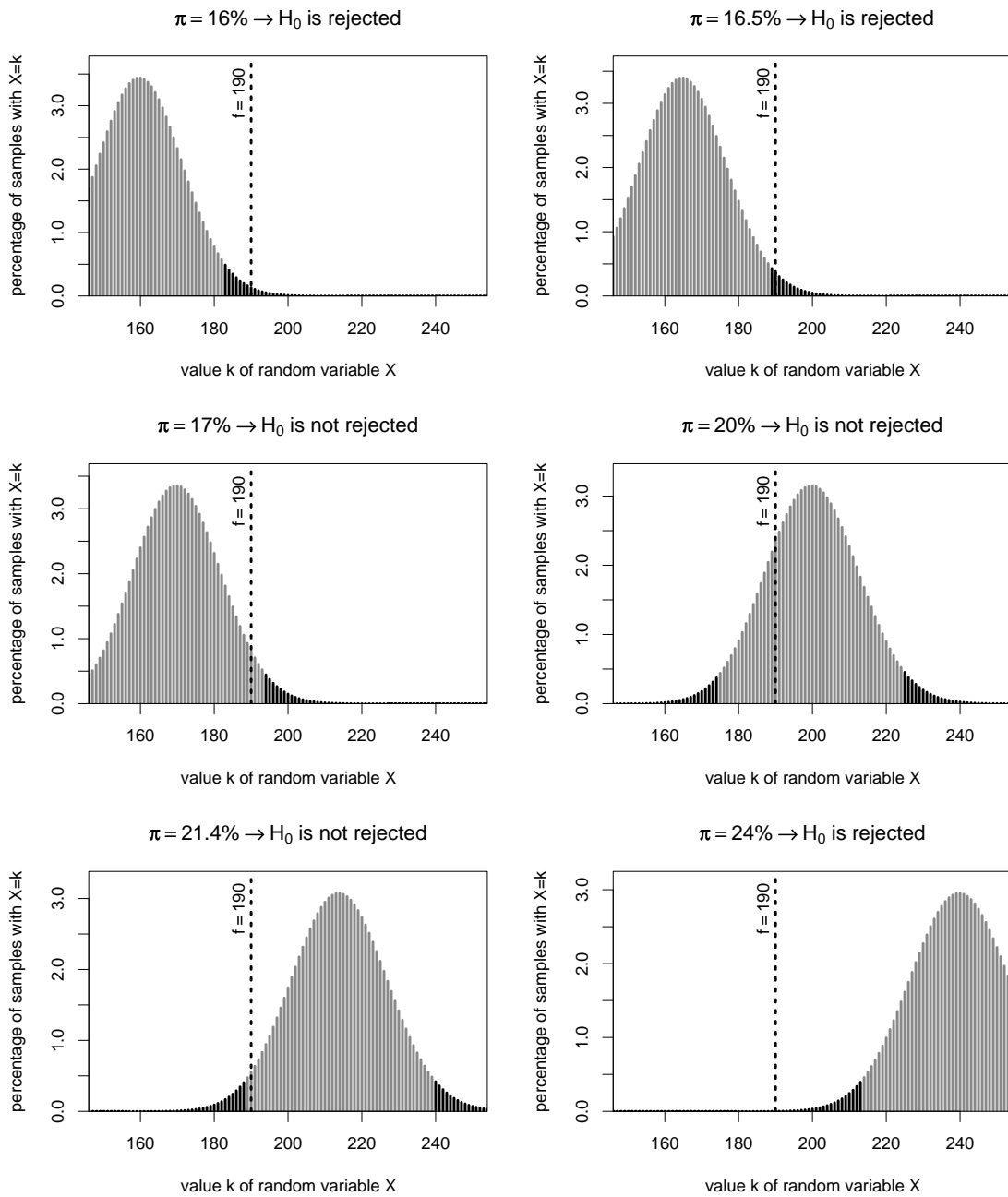


Figure 2: Illustration of the procedure for estimating a confidence set

	$n = 100$ $k = 19$	$n = 1,000$ $k = 190$	$n = 10,000$ $k = 1,900$
$\alpha = .05$	11.8% ... 28.1%	16.6% ... 21.6%	18.2% ... 19.8%
$\alpha = .01$	10.1% ... 31.0%	15.9% ... 22.4%	18.0% ... 20.0%
$\alpha = .001$	8.3% ... 34.5%	15.1% ... 23.4%	17.7% ... 20.3%

Table 1: Binomial confidence intervals for various sample sizes n and confidence levels α . The maximum-likelihood estimate is $\hat{\pi} = 19\%$ in each case.

1 test, this confidence set is an uninterrupted range of numbers and is called a *binomial*
2 *confidence interval*. Of course, it is infeasible to perform separate hypothesis tests for the
3 infinite number of possible null hypotheses $\pi = x$, but specialized mathematical algorithms
4 (available in all standard statistical software packages) can be used to compute the end
5 points of binomial confidence intervals efficiently. In our example, the observed data
6 $f = 190$ and $n = 1,000$ yield a confidence interval of $\pi \approx 16.6\% \dots 21.6\%$ (the common
7 mathematical notation for such a range, which you may encounter in technical literature,
8 is $[.166, .216]$).

9 The width of a binomial confidence interval depends on the sample size n and the
10 significance level α used in the test. As we have seen in Section 2, a larger value of
11 n makes it easier to reject the null hypothesis. Obviously, adopting a higher (i.e., less
12 conservative) value of α also makes it easier to reject H_0 . Hence these factors lead to
13 a narrower confidence interval (which, to reiterate this important point, consists of all
14 estimates x for which H_0 is *not* rejected). Table 1 shows confidence intervals for several
15 different sample sizes and significance levels. A confidence interval for a significance level of
16 $\alpha = .05$ (which keeps the risk of false rejection below 5%) is often called a 95% confidence
17 interval, indicating that we are 95% certain that the true population value π is somewhere
18 within the range (since we can rule out any other value with 95% certainty). Similarly, a
19 significance level of $\alpha = .01$ leads to a 99% confidence interval.

20 Confidence intervals can be seen as an extension of hypothesis tests. The 95% confi-
21 dence interval for the observed data immediately tells us whether a given null hypothesis
22 $H_0 : \pi = x$ would be rejected by the binomial test at significance level $\alpha = .05$. Namely,
23 H_0 is rejected if and only if the hypothesized value x does *not* fall within the confidence
24 interval. The width of a confidence interval illustrates thus how easily a null hypothesis
25 can be rejected, i.e., it gives an indication of how much the (unknown) true population
26 proportion π must differ from the value stipulated by the null hypothesis (which is often
27 denoted by the symbol π_0) so that H_0 will reliably be rejected by the hypothesis test.
28 Intuitively speaking, the difference between π and π_0 has to be considerably larger than
29 the width of one side of the 95% confidence interval so that it can reliably be detected by a
30 binomial test with $\alpha = .05$ (keep in mind that, even when the difference between π and π_0
31 is larger than this width, because of sampling variation, $\hat{\pi}$ and π_0 might be considerably
32 closer, leading to failure to reject H_0). The term *effect size* is sometimes used as a generic
33 way to refer to the difference between null hypothesis and true proportion. The reliability
34 of rejection given a certain effect and sample size is called the *power* of the hypothesis
35 test (see DeGroot/Schervish 2002, Chapter 8). In our example, the arithmetic difference
36 $\pi - \pi_0$ is a sensible way of quantifying effect size, but many other measures exist and may
37 be more suitable in certain situations (we will return to this issue during the discussion of
38 two-sample tests in Section 5).

39 In corpus analysis, we often deal with very large samples, for which confidence intervals
40 will be extremely narrow, so that a very small effect size may lead to highly significant

1 rejection of H_0 . Consider the following example: Baayen (2001, 163) claims that the
 2 definite article *the* accounts for approx. 6% of all words in (British) English, including
 3 punctuation and numbers. Verifying this claim on the LOB (the British equivalent of the
 4 Brown corpus, see article 22), we find highly significant evidence against H_0 . In particular,
 5 there are $f = 68,184$ instances of *the* in a sample of $n = 1,149,864$ words. A two-sided
 6 binomial test for $H_0 : \pi = 6\%$ rejects the null hypothesis with a p-value of $p \approx 0.1\%$.

7 However, the MLE for the true proportion π is actually very close to 6%, viz. $\hat{\pi} =$
 8 5.93%, and the 95% confidence interval is $\pi = 5.89\% \dots 5.97\%$. This difference is certainly
 9 not of scientific relevance, and $\hat{\pi}$ as well as the entire confidence range would be understood
 10 to fall under Baayen’s claim of “approximately 6%”. The highly significant rejection
 11 is merely a consequence of the large sample size and the corresponding high power of
 12 the binomial test. Gries (2005) is a recent discussion of the “significance” of statistical
 13 significance in corpus work.

14 At the opposite end of the scale, it is sometimes important to keep the sample size as
 15 small as possible, especially when the preparation of the sample involves time-consuming
 16 manual data annotation. Power calculations, which are provided by many statistical soft-
 17 ware packages, can be used to predict the minimum sample size necessary for a reliable
 18 rejection of H_0 , based on our conjectures about the true effect size.

19 4 The normal approximation

20 Looking back at Figure 1, we can see that the binomial sampling distribution has a fairly
 21 simple and symmetric shape, somewhat reminiscent of the outline of a bell. The peak of
 22 the curve appears to be located at the expected frequency $e = 15$. For other parameter
 23 values π and n , we observe the same general shape, only stretched and/or translated. This
 24 bell-shaped curve can be described by the following mathematical function:

$$f(x) := \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right) \quad (4)$$

25 This is the formula of a *normal* or *Gaussian distribution* (DeGroot/Schervish 2002,
 26 Section 5.6). The parameter μ , called the *mean*, will determine the peak of the bell-shaped
 27 curve, and the parameter σ , called the *standard deviation*, will determine the width of the
 28 curve (the symbol π in this formula stands for Archimedes’ constant $\pi = 3.14159\dots$ and
 29 not for a population proportion; to avoid another ambiguity, we write $\exp(-\frac{(x-\mu)^2}{2\sigma^2})$ for
 30 the exponential function in lieu of the more commonly encountered $e^{-(x-\mu)^2/2\sigma^2}$, since
 31 we are using e to denote the expected frequency). The roles of the two parameters are
 32 illustrated in Figure 3.

33 A binomial distribution with parameters n and π is approximated by a normal distri-
 34 bution with parameters $\mu = n\pi$ and $\sigma = \sqrt{n\pi(1-\pi)}$. Figure 4 shows the same binomial
 35 distribution illustrated in Figure 1 (with sample size $n = 100$ and proportion $\pi = 15\%$)
 36 and the corresponding normal approximation with parameters $\mu = 15$ and $\sigma \approx 3.57$. The
 37 quality of the approximation will increase with sample size and it will depend on π not
 38 being too skewed (i.e., not too close to 0 or 1). A rule of thumb might be to trust the
 39 approximation only if $\sigma > 3$, which is the case in our example (if you refer back to the
 40 formula for σ , you will notice that it depends, indeed, on n and the skewness of π).

41 The parameters of the normal approximation can be interpreted in an intuitive manner:
 42 μ coincides with the expected frequency e under H_0 (remember from Section 2 that e is also
 43 given by $n\pi$; we will use μ when referring to the normal distribution formula, e otherwise,
 44 but keep in mind that the two symbols denote the same quantity). σ tells us how much

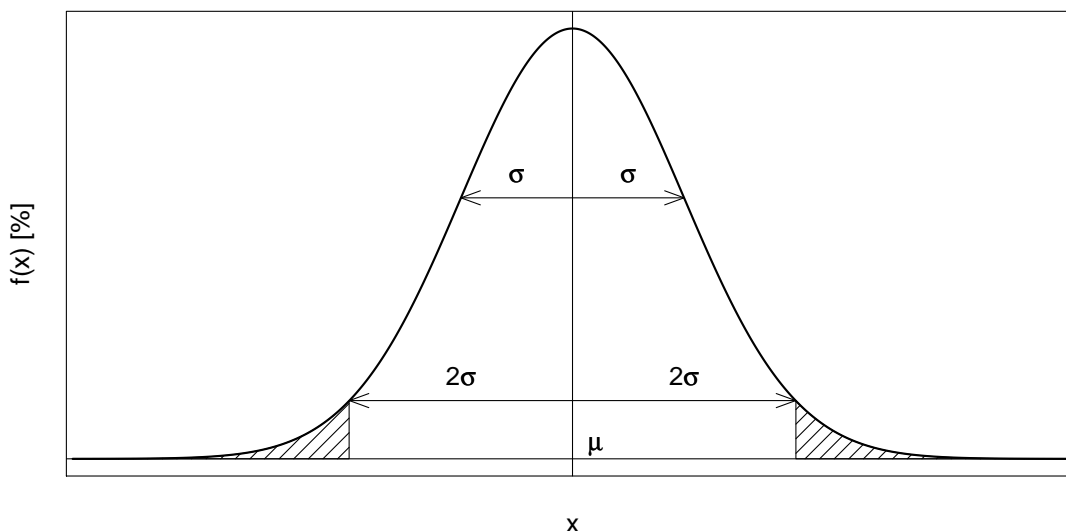


Figure 3: Interpretation of the parameters μ and σ of the normal distribution.

1 random variation we have to expect between different samples. Most of the samples will
 2 lead to observed frequencies between $\mu - 2\sigma$ and $\mu + 2\sigma$ and virtually all observed values
 3 will lie between $\mu - 3\sigma$ and $\mu + 3\sigma$ (refer to Figure 3 again), provided that H_0 is true.

4 To compute binomial tail probabilities based on a normal approximation, one calculates
 5 the corresponding area under the bell curve, as illustrated in Figure 4 for the tail proba-
 6 bility $\Pr(X \geq 19)$. In this illustration, we have also applied *Yates' continuity correction*
 7 (DeGroot/Schervish 2002, Section 5.8), which many statistical software packages use to
 8 make adjustments for the discrepancies between the smooth normal curve and the discrete
 9 distribution that is approximated. In our example, Yates' correction calculates the area
 10 under the normal curve for $x \geq 18.5$ rather than $x \geq 19$.

11 We find that the normal approximation gives a one-tailed p-value of 16.3% for observed
 12 frequency $f = 19$, sample size $n = 100$ and null hypothesis $H_0 : \pi = 15\%$. This is
 13 the same p-value we obtained from the (one-tailed) binomial test, indicating that the
 14 approximation is very good. Given that the normal distribution (unlike the binomial!)
 15 is always symmetrical, the two-tailed p-value can be obtained by simply multiplying the
 16 one-tailed value by two (which corresponds to adding up the tail areas under the curve for
 17 values that are at least as extreme as the observed value, with respect to e). In our case
 18 this gives 32.6%, again equivalent to the binomial test result.

19 There are two main reasons why the normal approximation is often used in place of the
 20 binomial test. First, the exact (non-approximated) binomial test and binomial confidence
 21 intervals require computationally expensive procedures that, for large sample sizes such
 22 as those often encountered in corpus-based work, can be problematic even for modern
 23 computing resources (a particularly difficult case is the extension of confidence intervals
 24 to the two-sample setting that we introduce in Section 5 and beyond). Second, the normal
 25 approximation leads to a more intuitive interpretation of the difference $f - e$ between
 26 observed and expected frequency, and the amount of evidence against H_0 that it provides
 27 (the importance of a given raw difference value depends crucially on sample size and on
 28 the null hypothesis proportion π_0 , which makes it hard to compare across samples and
 29 experiments).

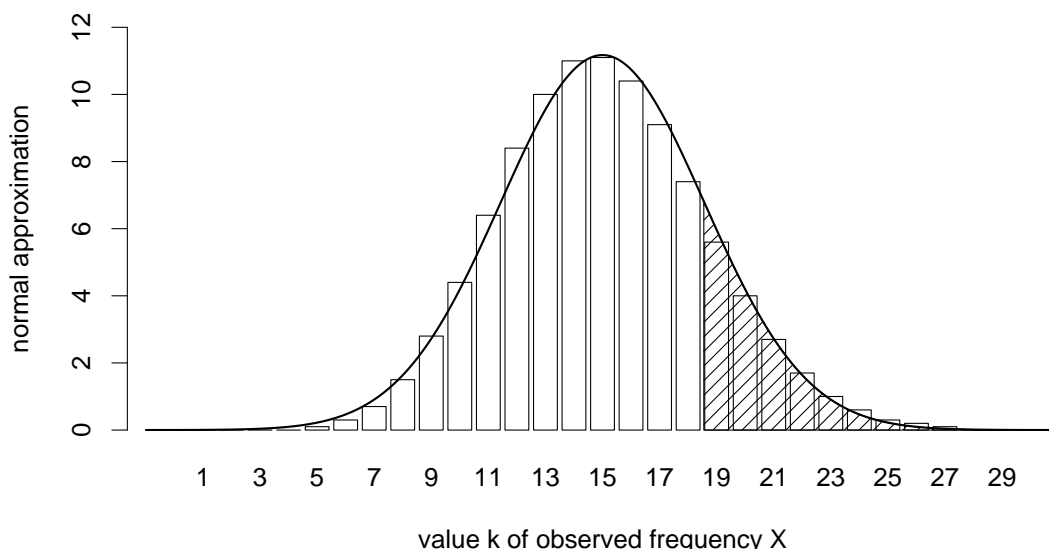


Figure 4: Approximation of binomial sampling distribution by normal distribution curve.

1 An interpretation of $f - e$ (or, equivalently $f - \mu$) that is comparable, e.g., between
 2 samples of different sizes, is achieved by a normalized value, the *z-score*, which divides
 3 $f - \mu$ by the standard deviation σ (you can think of this as expressing $f - \mu$ in σ 's, i.e.,
 4 using σ as the “unit of measurement”):

$$z := \frac{f - \mu}{\sigma} \quad (5)$$

5 If two observations f_1 and f_2 (possibly coming from samples of different sizes and
 6 compared against different null hypotheses) lead to the same z-score $z_1 = z_2$, they are
 7 equally “extreme” in the sense that they provide the same amount of evidence against their
 8 respective null hypothesis (as given by the approximate p-values). To get a feel for this,
 9 refer back to Figure 3, which illustrates the approximate two-tailed p-value corresponding
 10 to $z = 2$ as a shaded area under the normal curve. This area has exactly the same size
 11 regardless of the specific shape of the curve implied by H_0 (in the form of the parameters
 12 μ and σ). In other words, whenever we observe a value that translates into a z-score of
 13 $z = 2$ (according to the respective null hypothesis), we will obtain the same p-value from
 14 (the normal approximation to) the binomial test. Since we apply a two-tailed test, an
 15 observation that is two standard deviations to the left of the expected value ($z = -2$) will
 16 also lead to the same p-value.

17 Once an observation f has been converted into a z-score z , it is thus easy to decide
 18 whether H_0 can be rejected or not, by comparing $|z|$ with previously established *thresholds*
 19 for common significance levels α . For $\alpha = .05$, the (two-tailed) z-score threshold is 1.96, so
 20 the rejection criterion is $|z| \geq 1.96$; for $\alpha = .01$ the threshold is $|z| \geq 2.58$ and for $\alpha = .001$
 21 it is $|z| \geq 3.29$. Thus, no matter what the original values of f , π and n are, if in an
 22 experiment we obtain a z-score of, say, $z = 2$ (meaning that f is two standard deviations
 23 away from e), we immediately know that the result is significant at the .05 significance
 24 level, but not at the .01 level. Statistics textbooks traditionally provide lists of z-score
 25 thresholds corresponding to various significance levels, although nowadays p-values for
 26 arbitrary z-scores can quickly be obtained from statistical software packages.

1 5 Two-sample tests

2 Until now, we analyzed what is known as a *one-sample* statistical setting, where our null
3 hypothesis concerns a certain quantity (often, the proportion of a certain phenomenon)
4 in the set of all relevant units (e.g., all the sentences of English) and we use a sample of
5 such units to see if the null hypothesis should be rejected. However, *two-sample* settings,
6 where we have two samples (e.g., two corpora with different characteristics, or two sets
7 of sentences of different types) and want to know whether they are significantly different
8 with respect to a certain property, are much more common. Coming back to the example
9 of passivization in idiomatic vs. non-idiomatic constructions from the introduction, our
10 two samples would be sets of idiomatic and non-idiomatic constructions; we would count
11 the number of passives in both sets; and we would verify the null hypothesis that there is
12 no difference between the proportion of passives in the two samples.

13 It is easier to motivate one aspect of hypothesis testing that is often counter-intuitive,
14 i.e., the fact that we pick as null hypothesis the “uninteresting” hypothesis that we hope
15 to reject, when looking at the two-sample case. First, most linguistic theories, especially
16 categorical ones, are more likely to predict that there is *some* difference between two sets,
17 rather than making quantitative predictions about this difference being of a certain size.
18 Second, in this way, if we can reject the null hypothesis, we can claim that the hypothesis
19 that there is no difference between the groups is not tenable, i.e., that there *is* a difference
20 between the groups, which is what our theory predicts. If, instead, we tested the null
21 hypothesis that there is a certain difference between the groups, and we found that this
22 hypothesis cannot be rejected, we could only claim that, for now, we have not found
23 evidence that would lead us to reject our hypothesis: clearly, a weaker conclusion.

24 Probably the majority of questions that are of interest to linguists can be framed
25 in terms of a two-sample statistical test: for several examples of applications in syntax,
26 see article 45; for an application to the study of collocations, see article 57. Here, we
27 discuss the example of the distribution of passives in two broad classes of written English,
28 “informative” prose (such as daily press) and “imaginative” prose (such as fiction). One
29 plausible *a priori* hypothesis is that these two macro-genres will differ in passive ratios,
30 with a stronger tendency to use passives in informative prose, due to the impersonal, more
31 “objective” tone conferred to events by passive voice (for a more serious corpus-based
32 account of the distribution of the English passive, including register-based variation, see
33 Biber/Johansson/Leech/et al. 1999, Sections 6.4 and 11.3). Our null hypothesis will be
34 that there is no difference between the proportion of passives in informative prose π_1 and
35 imaginative prose π_2 , i.e., $H_0 : \pi_1 = \pi_2$. Conveniently, the documents in the Brown corpus
36 are categorized into informative and imaginative writing – thus, we can draw random
37 samples of $n_1 = 100$ sentences from the informative section of the corpus, and $n_2 = 100$
38 sentences from the imaginative section. Counting the passives, we find that the informative
39 sample contains $f_1 = 23$ passives, whereas the imaginative sample contains $f_2 = 9$ passives.

40 Since there is a considerable difference between f_1 and f_2 , we are tempted to reject
41 H_0 . However, before we can do so, we must find out to what extent the difference can
42 be explained by random variation, i.e., we have to calculate how likely it is that the two
43 samples come from populations with the same proportion of passives, as stated by the
44 null hypothesis (statistics textbooks will often phrase the null hypothesis directly as: the
45 samples *are* from the same population). In order to calculate expected frequencies, we
46 have to estimate this common value from the available data, using maximum-likelihood
47 estimation: $\hat{\pi} = (f_1 + f_2)/(n_1 + n_2) = 32/200 = 16\%$ (we sum the f 's and n 's because,
48 if H_0 is right, then we can treat all the data we have as a larger sample from what, for
49 our purposes, counts as the same population). Replacing H_0 by the more specific null

1 hypothesis $H'_0 : \pi_1 = \pi_2 = 16\%$, we can compute the expected frequencies under the null
 2 hypothesis, i.e., $e_1 = e_2 = 100 \cdot \hat{\pi} = 16$ (which are identical in our case since $n_1 = n_2 = 100$),
 3 as well as the binomial sampling distributions.

4 In the one-sample case, we looked at the overall probability of f and all other possible
 5 values that are more extreme than $|f - e|$. The natural extension to the two-sample case
 6 would be to look at the overall probability of the pair (f_1, f_2) and all the other possible pairs
 7 of values that, taken together, are more extreme than the sum of $|f_1 - e_1|$ and $|f_2 - e_2|$.
 8 The lower this probability, the more confident we can be that the null hypothesis is false.
 9 In our case, $|f_1 - e_1|$ and $|f_2 - e_2|$ are directly comparable and might be added up in this
 10 way, since the expected frequencies $e_1 = e_2 = 16$ and the sample sizes $n_1 = n_2 = 100$
 11 are the same. However, in many real life situations, we will have to deal with samples of
 12 (sometimes vastly) different sizes (e.g., if one of the conditions is relatively rare so that
 13 only few examples can be found).

14 Fortunately, we know a solution to this problem from Section 4: z-scores provide
 15 a measure of extremeness that is comparable between samples of different sizes. We
 16 thus compute the z-scores $z_1 = (f_1 - e_1)/\sigma_1$ and $z_2 = (f_2 - e_2)/\sigma_2$ (with σ_1 and σ_2
 17 obtained from the estimate $\hat{\pi}$ according to H'_0). For mathematical reasons, the total
 18 extremeness is computed by adding up the squared z-scores $x^2 := (z_1)^2 + (z_2)^2$ instead of
 19 the absolute values $|z_1|$ and $|z_2|$. It should be clear that the larger this value is, the less
 20 likely the null hypothesis of no difference in population proportions is, and thus we should
 21 feel more confident in rejecting it. More precisely, the p-value associated with x^2 is the
 22 sum over the probabilities of all outcomes for which the corresponding random variable
 23 $X^2 := (Z_1)^2 + (Z_2)^2$ is at least as large as the observed x^2 , i.e. $\Pr(X^2 \geq x^2)$.

24 Instead of enumerating all possible pairs of outcomes with this property, we can again
 25 make use of the normal approximation, which leads to a so-called *chi-squared distribution*
 26 with one *degree of freedom* ($df = 1$). Using the chi-squared distribution, we can easily
 27 calculate the p-value corresponding to the observed x^2 , or compare x^2 with known rejection
 28 thresholds for different significance levels (e.g. $x^2 \geq 3.84$ for $\alpha = .05$ or $x^2 \geq 6.63$ for
 29 $\alpha = .01$). This procedure is known as *(Pearson's) chi-squared test* (Agresti 1996, Section
 30 2.4; DeGroot/Schervish 2002, Sections 9.1-4).

An alternative representation of the observed frequency data that is widely used in
 statistics takes the form of a so-called *contingency table*:

	sample 1	sample 2	
passives	f_1	f_2	(6)
other	$n_1 - f_1$	$n_2 - f_2$	

31 The cells in the first row give the frequencies of passives in the two samples, while the
 32 cells in the second row give the frequencies of all other sentence types. Notice that each
 33 column of the contingency table adds up to the respective sample size, and that $\hat{\pi}$ (the
 34 estimated population proportion under H_0 needed to compute expected frequencies) can
 35 be obtained by summing over the first row and dividing by the overall total. Thus, the
 36 chi-squared statistic x^2 can easily be calculated from such a table (Agresti 1996, Chapter
 37 2; DeGroot/Schervish 2002, Section 9.3) and most statistical software packages expect
 38 frequency data for the chi-squared test in this form. Like in the one-sample case, the
 39 normal approximation is only valid if the sample sizes are sufficiently large. The standard
 40 rule of thumb for contingency tables is that all *expected* cell frequencies (under H'_0) must be
 41 ≥ 5 (Agresti 1996, Section 2.4.1). In the usual situation in which $\hat{\pi} < 50\%$, this amounts
 42 to $n_1 \hat{\pi} \geq 5$ and $n_2 \hat{\pi} \geq 5$. Statistical software will usually produce a warning when the
 43 normal approximation is likely to be inaccurate.

1 There is also an *exact* test for contingency tables, similar to the binomial test in the
 2 one-sample case. This test is known as *Fisher's exact test* (Agresti 1996, Section 2.6). It
 3 is implemented in most statistical software packages, but it is computationally expensive
 4 and may be inaccurate for large samples (depending on the specific implementation).
 5 Therefore, use of Fisher's test is usually reserved for situations where the samples are too
 6 small to allow the normal approximations underlying the chi-squared test (as indicated by
 7 the rule of thumb above).

In the current example ($f_1 = 23$ and $f_2 = 9$), the contingency table corresponding to the observed data is

	sample 1	sample 2	
passives	23	9	(7)
other	77	91	

8 Using a statistical software package, we obtain $x^2 = 6.29$ for this contingency table,
 9 leading to rejection of H_0 at the .05 significance level (but not at the .01 level). The
 10 approximate p-value computed from x^2 is $p = 1.22\%$, while Fisher's exact test yields
 11 $p = 1.13\%$ (with expected frequencies $n_1\hat{\pi} = n_2\hat{\pi} = 16 \gg 5$, we anticipated a good
 12 agreement between the exact and the approximate test). We can thus conclude that there
 13 is, indeed, a difference between the proportion of passives in informative vs. imaginative
 14 prose. Moreover, the direction of the difference confirms our conjecture that the proportion
 15 of passives is higher in informative prose.

16 A particular advantage of the contingency table notation is that it allows straight-
 17 forward generalizations of the two-sample frequency comparison. One extension is the
 18 comparison of more than two samples representing different conditions (leading to a con-
 19 tingency table with $k > 2$ columns). For instance, we might want to compare the frequency
 20 of passives in samples from the six subtypes of imaginative prose in the Brown corpus (gen-
 21 eral fiction, mystery, science fiction, etc.). The null hypothesis for such a test is that the
 22 proportion of passives is the same for all six subtypes, i.e. $H_0 : \pi_1 = \pi_2 = \dots = \pi_6$.
 23 Another extension leads to contingency tables with $m > 2$ rows. In our example, we have
 24 distinguished between passive sentences on one hand and all other types of sentences on
 25 the other. However, this second group is less homogeneous so that further distinctions may
 26 be justified, e.g., at least between sentences with intransitive and transitive constructions.
 27 From such a three-way classification, we would obtain three frequencies $f^{(p)}$, $f^{(i)}$ and $f^{(t)}$
 28 for each sample, which add up to the sample size n . These frequencies can naturally be
 29 collected in a contingency table with three rows. The null hypothesis would now stipu-
 30 late that the proportions of passives, intransitive and transitives are the same under both
 31 conditions (assuming $k = 2$), viz. $\pi_1^{(p)} = \pi_2^{(p)}$, $\pi_1^{(i)} = \pi_2^{(i)}$ and $\pi_1^{(t)} = \pi_2^{(t)}$. In general, an
 32 x^2 value can be calculated for any $m \times k$ contingency table in analogy to the 2×2 case.
 33 The p-value corresponding to x^2 can be obtained from a chi-squared distribution with
 34 $df = (m - 1)(k - 1)$ degrees of freedom. If the expected frequency in at least one of the
 35 cells is less than 5, a version of Fisher's exact test can be used (this version is considerably
 36 *more* expensive than Fisher's test for 2×2 tables, though).

37 Having established that the proportion of passives is different in informative vs. imagi-
 38 native prose, we would again like to know how large the effect size is, i.e. by how much the
 39 proportions π_1 and π_2 differ. This is particularly important for large samples, where small
 40 (and hence linguistically irrelevant) effect sizes can easily lead to rejection of H_0 (cf. the
 41 discussion in Section 3). A straightforward and intuitive measure of effect size is the dif-
 42 ference $\delta := \pi_1 - \pi_2$. When the sample sizes are sufficiently large, normal approximations
 43 can be used to compute a confidence interval for δ . This procedure is often referred to as
 44 a *proportions test* and it is illustrated, for example, by Agresti (1996, Section 2.2). In our

1 example, the 95% confidence interval is $\delta = 3.0\% \dots 25.0\%$, showing that the proportion
2 of passives is at least 3 percentage points higher in informative prose than in imaginative
3 prose (with 95% certainty).

4 In other situations, especially when π_1 and π_2 are on different orders of magnitude,
5 other measures of effect size, such as the ratio π_1/π_2 (known as *relative risk*) may be more
6 appropriate. A related measure, the *odds ratio* θ , figures prominently because an exact
7 confidence interval for θ can be obtained from Fisher's test. Most software packages that
8 implement Fisher's test will also offer calculation of this confidence interval. In many
9 linguistic applications (where π_1 and π_2 are relatively small), θ can simply be interpreted
10 as an approximation to the ratio of proportions (relative risk), i.e., $\theta \approx \pi_1/\pi_2$. On these
11 measures see, again, Section 2.2 of Agresti (1996). Effect size in general $m \times k$ contingency
12 tables is much more difficult to define, and it is most often discussed in the setting of
13 so-called *generalized linear models* (Agresti 1996, Chapter 4).

14 Examples of fully worked two-sample analyses based on contingency tables can be
15 found in articles 45 and 57. As illustrated by article 45 in particular, contingency tables
16 and related two-sample tests can be tuned to a number of linguistic questions by look-
17 ing at different kinds of linguistic populations. For example, if we wanted to study the
18 distribution of by-phrases in passive sentences containing two classes of verbs (say, verbs
19 with an agent vs. experiencer external argument), we could define our two populations
20 as all passive sentences with verbs of class 1 and all passive sentences with verbs of class
21 2. We would then sample passive sentences of these two types, and count the number of
22 by-phrases in them. As a further example, we might be interested in comparing alterna-
23 tive morphological and syntactic means to express the same meaning. For example, we
24 might be interested, with Rainer (2003), in whether various classes of Italian adjectives
25 are more likely to be intensified by the suffix *-issimo* or by the adverb *molto*. This leads
26 naturally to a contingency table for intensified adjectives with *-issimo* and *molto* columns,
27 and as many rows as the adjective classes we are considering (or vice versa). The key
28 to the successful application of statistical techniques to linguistic problems lies in being
29 able to frame interesting linguistic questions in operational terms that lead to meaning-
30 ful significance testing. The following section will discuss different ways to perform this
31 operationalization.

32 6 Linguistic units and populations

33 As we just said, from the point of view of linguists interested in analyzing their data
34 statistically, the most important issue is how to frame the problem at hand so that it can
35 be operationalized in terms suitable for a statistical test. In this section, we introduce
36 some concepts that might be useful when thinking of linguistic questions in a statistical
37 way.

38 In the example used throughout the preceding sections, we have defined the population
39 as the set of all (written American) English sentences and considered random samples of
40 sentences from this population. However, statistical inference can equally well be based
41 on any other linguistic unit, such as words, phrases, paragraphs, documents, etc. This
42 *unit of measurement* is often called a *token* in corpus linguistics, at least when referring to
43 words. Here, we use the term more in general to refer to any unit of interest.

44 The *population* then consists of all the utterances that have ever been produced (or
45 could be produced) in the relevant (sub)language, broken down into tokens of the chosen
46 type. We might also decide to focus on tokens that satisfy one or more other criteria
47 and narrow down the population to include only these tokens. For instance, we might be

1 concerned with the population of words that belong to a specific syntactic category; or
2 with sentences that contain a particular verb or construction, etc.

3 What we are interested in is the *proportion* π of tokens (in the population) that have
4 a certain additional property: e.g., word tokens that are nouns, verb tokens that belong
5 to the inflectional paradigm of *to make*, sentences in the passive voice, etc. The properties
6 used to categorize tokens for this purpose are referred to as *types* (in contrast to tokens,
7 which are the categorized objects).

8 Since the full population is inaccessible, our conclusions have to be based on a (*random*)
9 *sample* of tokens from the population. Such a sample of language data is usually called a
10 *corpus* (or can be derived from a corpus: when we define the population as a set of verb
11 tokens, for example, our sample might comprise all instances of verbs found in the corpus).
12 The *sample size* n is the total number of tokens in the sample, and the number of tokens
13 that exhibit the property of interest (i.e., that belong to the relevant type) is the observed
14 *frequency* f .

15 The same observed frequency can have different interpretations (with respect to the
16 corresponding population proportion) depending on the units of measurement chosen as
17 tokens, and the related target population. For instance, the number of passives in a
18 sample could be seen relative to the number of sentences (π = proportion of sentences
19 in the passive voice), relative to the number of verb phrases (π = proportion of passive
20 verb phrases), relative to word tokens (π = relative frequency of passive verb phrases per
21 1,000 words), relative to all sentences containing transitive verbs (π = relative frequency
22 of actual passives among sentences that could in principle be in passive voice). Note that
23 each of these interpretations casts a different light on the observed frequency data. It
24 is the linguist's task to decide which interpretation is the most meaningful, and to draw
25 conclusions about the linguistic research questions that motivated the corpus study.

26 Other examples might include counting the number of deverbal nouns in a sample
27 from the population of all nouns in a language; counting the number of words ending in
28 a coronal stop in a sample from the population of all words in the language; counting the
29 number of sentences with heavy-NP shift in a sample from the population of all sentences
30 with a complement that could in principle undergo the process; counting the number
31 of texts written in first person in a sample from the population of literary texts in a
32 certain language and from a certain period. Related problems can also be framed in terms
33 of looking at *two* samples from distinct populations (cf. Section 5), e.g., counting and
34 comparing the number of deverbal nouns in samples from the populations of abstract and
35 concrete nouns; counting the number of words ending in a coronal stop in samples from the
36 population of all native words in the language and the population of loanwords; counting
37 the number of texts written in first person in samples from populations of texts belonging
38 to two different literary genres.

39 In many cases, frequencies are computed not only for a single property, but for a set
40 of mutually exclusive properties, i.e., a *classification* of the tokens into different types. In
41 the two-sample setting this leads naturally to a $m \times 2$ contingency table (with the types
42 in the classification as rows, and the two populations we are comparing as columns). Note
43 that the classification has to be *complete*, so that the columns of the table add up to the
44 respective sample sizes, which is often achieved by introducing a category labeled "other"
45 (the single-property/two-samples cases above correspond to 2×2 contingency tables with
46 an "other" class: e.g., deverbal vs. "other" nouns compared across the populations of
47 abstract vs. concrete nouns).

48 As an example of a classification into multiple categories, word tokens might be clas-
49 sified into syntactic categories such as noun, verb, adjective, adverb, etc., with an "other"

1 class for minor syntactic categories and problematic tokens. A chi-squared test might then
2 be performed to compare the frequencies of these categories in samples from two genres.
3 As another example, one might classify sentences according to the semantic class of their
4 subject, and then compare the frequency of these semantic classes in samples of the pop-
5 ulations of sentences headed by true intransitive vs. unaccusative verbs. It is not always
6 obvious which characteristics should be operationalized as a classification of the tokens
7 into types, and which should rather be operationalized in terms of different populations
8 the tokens belong to. In some cases, it might make more sense to frame the task we just
9 discussed in terms of the distribution of verb types across populations of sentences with
10 different kinds of subjects, rather than vice versa. This decision, again, will depend on the
11 linguistic question we want to answer.

12 In corpus linguistics, *lexical classifications* also play an important role. In this case,
13 types are the distinct word forms or lemmas found in a corpus (or sequences of word forms
14 or lemmas). Lexical classifications may lead to extremely small proportions π (sometimes
15 measured in occurrences per million words) and huge differences between populations in
16 the two-sample setting. Article 57 discusses some of the relevant methodologies in the
17 context of collocation extraction.

18 The examples we just discussed give an idea of the range of linguistic problems that
19 can be studied using the simple methods based on count data described in this article.
20 Other problems (or the same problems viewed from a different angle) might require other
21 techniques, such as those mentioned in the next two sections. For example, our study of
22 passives could proceed with a *logistic regression* (see Section 8), where we look at which
23 factors have a significant effect on whether a sentence is in the passive voice or not. In any
24 case, it will be fundamental for linguists interested in statistical methods to frame their
25 questions in terms of populations, samples, types and tokens.

26 7 Non-randomness and the unit of sampling

27 So far, we have always made the (often tacit) assumption that the observed data (i.e., the
28 corpus) are a random sample of tokens of the relevant kind (e.g., in our running example
29 of passives, a sentence) from the population. Most obviously, we have compared a corpus
30 study to drawing balls from an urn in Section 2, which allowed us to predict the sampling
31 distribution of observed frequencies. However, a realistic corpus will rarely be built by
32 sampling individual tokens, but rather as a collection of contiguous stretches of text or
33 even entire documents (such as books, newspaper editions, etc.). For example, the Brown
34 corpus consists of 2,000-word excerpts from 500 different books (we will refer to these
35 excerpts as “texts” in the following). The discrepancy between the *unit of measurement*
36 (a token) and the *unit of sampling* (which will often contain hundreds or thousands of
37 tokens) is particularly obvious for lexical phenomena, where tokens correspond to single
38 words. Imagine the cost of building the Brown corpus by sampling a single word each
39 from a million different books rather than 2,000 words each from only 500 different books!

40 Even in our example, where each token corresponds to an entire sentence, the unit of
41 sampling is much larger than the unit of measurement: each text in the Brown contains
42 roughly between 50 and 200 sentences. This need not be a problem for the statistical
43 analysis, as long as each text is itself a random sample of tokens from the population, or
44 at least sufficiently similar to one. However, various factors, such as the personal style of
45 an author or minor differences in register or conventions within a particular subdomain,
46 may have a *systematic* influence on how often passives are used in different texts. This
47 means that the variability of the frequency of passives between texts may be much larger

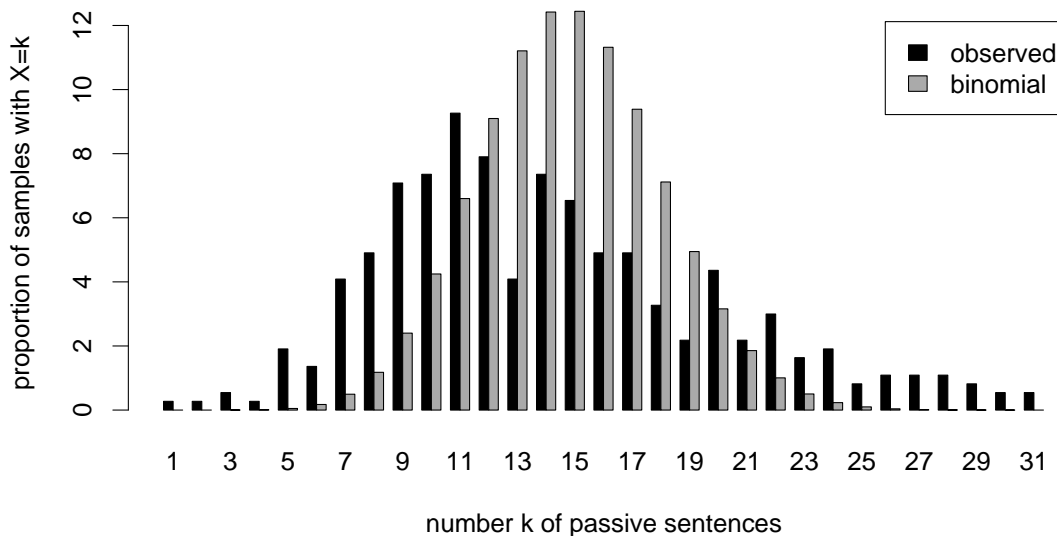


Figure 5: Comparison of the frequencies of passives in the texts of the Brown corpus (informative prose only) with the binomial distribution predicted for random samples. In order to ensure comparability of the frequencies, 50 sentences were sampled from each Brown text.

1 than between random samples of the same sizes (where all variation is purely due to chance
 2 effects).

3 Again, the problem is most obvious for lexical frequencies. Many content words (except
 4 for the most general and frequent ones) will almost only be found in texts that deal with
 5 suitable topics (think of nouns like *football* or *sushi*, or adjectives like *coronal*). On the
 6 other hand, such topical words tend to have multiple occurrences in the same text, even if
 7 these would be extremely unlikely in a random sample (indeed, the “burstiness” of words
 8 in specific texts is used as strategy to find interesting keywords; see, e.g., Church 2000).

9 The increased variability of frequency between the individual texts is attenuated to
 10 some extent when corpus frequencies are obtained by summing over all the (randomly
 11 selected) texts in a corpus. However, in most cases the corpus frequencies will still show
 12 more variation than predicted by the binomial distribution.

13 In order to verify empirically whether a linguistic phenomenon such as the frequency of
 14 passives is subject to such non-randomness, we can compare the distribution of observed
 15 frequencies across the texts in a corpus with the distribution predicted for random samples
 16 by the binomial distribution. An example of such a comparison is shown in Figure 5.
 17 From each Brown text, we have taken a sample of 50 sentences (this subsampling step
 18 was necessary because the number of sentences per text varies from around 50 to more
 19 than 200). By tabulating the observed frequencies, we obtain the distribution shown as
 20 black bars in Figure 5. The gray bars show the binomial distribution that we would
 21 have obtained for random samples from the full population (the population proportion of
 22 passives was estimated at $\pi = 27.5\%$, based on the Brown data). Note that we have used
 23 only informative prose texts, since we already know from Section 5 that the proportion of
 24 passives differs considerably between the two major sections of the corpus.

25 As the figure shows, the observed amount of variation is larger than the one predicted
 26 from the binomial distribution: look for example at the proportion of observed and bino-

1 mial samples with $X \leq 7$. The standard deviation (which, as discussed in Section 4, is
2 a measure of the width of a distribution) is $\sigma = 6.63$ for the empirical distribution, but
3 only $\sigma = 3.16$ for the binomial distribution. The corresponding z-scores, having σ in the
4 denominator (see Equation (5) in Section 4), will be smaller for the empirical distribution,
5 and thus the results are less significant than they would seem according to the binomial dis-
6 tribution. This means that the binomial test will lead to rejection of a true null hypothesis
7 more easily than should be the case, given the spread of the actual distribution.

8 Suppose that we want to test the null hypothesis $H_0 : \pi = 27.5\%$ (which is in fact
9 true) based on a sample of $n = 50$ sentences from the informative prose in the Brown
10 corpus. If the observed frequency of passives in this sample is $f = 7$, we feel confident
11 to reject H_0 ($e = 13.75$ leads to a z-score of $z = -2.14$, above the $\alpha = .05$ threshold of
12 $|z| \geq 1.96$). However, if all sentences in this sample came from the same text (rather than
13 being sampled randomly from the entire informative prose section), Figure 5 shows that
14 the risk of obtaining *precisely* $f = 7$ by chance is already around 4%! The “true” z-score
15 (based on the standard deviation computed from the observed samples) is only $z = -1.02$,
16 far away from any rejection threshold (in fact, this z-score indicates a risk of more than
17 30% that H_0 would be wrongly rejected).

18 Seeing how non-randomness effects can lead to a drastic misinterpretation of the ob-
19 served frequency data, a question arises naturally: How can we make sure that a corpus
20 study is not affected by non-randomness? While for many practical purposes it might be
21 possible to ignore the issue, the only way to be absolutely sure is to ascertain that the unit
22 of sampling coincides with the unit of measurement. When using a pre-compiled corpus
23 (as will be the case for most studies in corpus linguistics) or when it would be prohibitively
24 difficult and time consuming to sample individual tokens, we have no choice but to adjust
25 the *unit of measurement*. For example, when our data are based on the Brown corpus, the
26 unit of sampling – and hence the unit of measurement – would be a text, i.e., a 2,000-word
27 excerpt from a coherent document. Of course, we can no longer classify such an excerpt
28 as “passive” or “non-passive”. Instead, what we observe for each token is a real-valued
29 number: the proportion of passive sentences in the text.

30 Unlike previously, where each measurement was essentially a yes/no-decision (“passive”
31 or “not passive”) or a m -way classification, measurements are now real-valued numbers
32 that can in principle assume any value between 0 and 1 ($3/407$, $139/211$, etc.). Statisticians
33 speak of a *nominal scale* (for yes/no-decisions and classifications) vs. an *interval scale* (for
34 numerical data). In order to analyze such data, we need an entirely different arsenal of
35 statistical procedures, such as the well-known *t-test*. These methods are explained in any
36 introductory textbook on statistics, and we give a brief overview of the most important
37 terms and techniques in Section 8.

38 This approach is only viable for phenomena, such as passive voice, that have a rea-
39 sonably large number of occurrences in each text. It would not be sensible to count the
40 proportion of occurrences of the collocation *strong tea* in the Brown texts (or even in a
41 corpus made of larger text stretches), since the vast majority of texts would yield a pro-
42 portion of 0% (in the Brown corpus, *strong tea* occurs exactly once, which means that in
43 all texts but one the proportion will indeed be 0%).

44 Notice that, from a statistical perspective, the issues of representativeness and balance
45 sometimes discussed in connection to corpus design (see article 11) involve two aspects: 1)
46 How to define the target population precisely (is it possible to delimit a set of utterances
47 that constitutes the population of, say, “contemporary English”?), and 2) how to take a
48 random sample from the target population (with the complication discussed in this section
49 that what might constitute a random sample of, say, documents, will not be a random

1 sample of, say, sentences). See Evert (2006) for an extended discussion of non-randomness
2 in corpus linguistics.

3 8 Other techniques

4 Like for count data, there is a range of statistical tests that can be used to analyze data on
5 an interval scale (such as the relative number of sentences containing passives per document
6 discussed in the previous section, or reaction times from a psycholinguistic experiment).
7 For the one-sample case, in which we want to test whether an observed interval-scale
8 quantity (such as the proportion of passive sentences in a text) could plausibly come from
9 a population where the same quantity has a certain distribution with a specific mean
10 and standard deviation, you can use a *(one sample) t-test* (comparable to the binomial
11 test for count data or, more precisely, the normal approximation based on z-scores). Un-
12 surprisingly, when two samples are compared, the appropriate test is a *two-sample t-test*
13 (corresponding to the chi-squared test for count data). However, in order to compare
14 more than two samples, rather than performing a series of pairwise t-tests (a procedure
15 that would make it much more likely that we obtain a significant result by chance), the
16 technique to be applied is the *one-way analysis of variance (ANOVA)*. The ANOVA can
17 only tell us whether at least one sample in the set is different from at least one other
18 sample, and post-hoc tests must then be performed to identify the sample(s) responsible
19 for rejection of H_0 .

20 In some settings, the variables in two samples have a natural pairing. For example, if
21 we compare the proportion of passives in English and Spanish texts based on a parallel
22 corpus, we should make use of the information that the texts are paired in order to control
23 for irrelevant factors that may affect passive proportion (e.g., style and topic of a text),
24 which should have a similar effect in an original and its translation. The appropriate test,
25 in this case, is the *paired t-test*.

26 In many studies, it makes sense to operationalize the problem as one of assessing the
27 association between two properties of the same unit, both measured on an interval scale.
28 For example, we might be interested in the issue of whether there is a relation between the
29 proportion of passives and, say, that of nominalizations (as both are plausible markers of
30 more formal registers). Given a list of pairs specifying the (relative) frequencies of passives
31 and nominalizations in each of the texts in our sample, we can perform a *correlation*
32 *analysis*. In this case, the null hypothesis will be that there is no correlation between the
33 two properties; and effect size will be measured in terms of how much of the variability
34 of one variable can be explained by linear dependence on the other variable (standard
35 correlation analysis will not capture *nonlinear* relations between variables).

36 A significant correlation does not imply a causal relation between two variables (even
37 if the numbers of passives and nominalizations turn out to be correlated, it is unlikely that
38 passives “cause” nominalization or vice versa). Often, however, we want to go beyond
39 the mere observation of a relationship between two variables. If we hypothesize that the
40 behavior of a certain variable depends on that of one or more other variables, we will
41 want to use statistics to test whether our *independent* variables predict the values of the
42 *dependent* variable beyond chance level. In this case, we use the technique of *(multiple)*
43 *linear regression* (which is related to correlation). In linear regression, the independent
44 variables can be a mixture of discrete and continuous variables, but the dependent variable
45 must be continuous.

46 Similar techniques can also be applied to the analysis of the kind of categorical data
47 (resulting in a contingency table of frequency counts) that have been the focus of this

1 article. The equivalent of linear regression in this case is *logistic regression*. For example,
2 a logistic regression analysis could try to predict whether a sentence is in passive voice or
3 not (a dichotomous dependent variable) in terms of factors such as the semantic class of the
4 verb (a categorical variable), the overall entropy of the sentence (a continuous variable),
5 etc. A full regression analysis tests significant effects of the independent variables, but
6 typically it also checks that the independent variables are not correlated with each other,
7 and it might look for the optimal combination of independent variables.

8 The cases we listed here (detection of differences and estimation of population values
9 in one/two/multiple paired/non-paired sample cases, assessment of association/correlation
10 between variables, regression) constitute a nearly exhaustive survey of the analytical set-
11 tings considered in inferential statistics. More advanced techniques, rather than introduc-
12 ing completely new scenarios, will typically deal with cases in which one or more of the as-
13 sumptions of the basic models are not met or the models need to be extended. For example,
14 more sophisticated ANOVA models can take multiple categorizations of the data and their
15 interactions into account (akin to the analysis of $m \times k$ contingency tables with m and/or
16 k greater than 2). Advanced regression techniques can detect non-linear relations between
17 the dependent and independent variables. So-called “distribution-free” tests make no as-
18 sumption about the distribution of the underlying population(s) nor about the sampling
19 distribution (these are typically referred to as *non-parametric methods*). Simulation-based
20 methods (*Monte Carlo methods*, the *Bootstrap*) provide an alternative to analytical esti-
21 mation of various parameters. A wealth of exploratory and visual methods are available
22 to evaluate the validity of assumptions and the quality of the resulting models. *Bayesian*
23 *inference*, a very important branch of statistics, allows, among other things, to distin-
24 guish between “more plausible” and “less plausible” estimates within a confidence interval
25 (the classic binomial confidence interval described in Section 3 indicates a range of plausi-
26 ble values for the population proportion, but it does not distinguish among these values,
27 whereas, intuitively, we would consider the MLE proportion much more plausible than,
28 say, the values at the edges of the confidence interval).

29 Some important kinds of corpus data, such as distributions of word types, are charac-
30 terized by the presence of a very large number of very rare types (words that occur only
31 once or never at all in the corpus at hand) and few extremely frequent types (function
32 words). These extremely skewed distributions make the application of standard statistical
33 models to certain tasks problematic (mainly, estimating the number of word types in a
34 population as well as related quantities), and demand specialized statistical tools. For a
35 general survey of the problems involved, see article 39 and the references on statistical
36 modeling of word frequency distributions recommended there.

37 Almost every elementary statistics textbook (including those listed in the next sec-
38 tion) will introduce t-tests, ANOVA, correlation and regression. Advanced techniques are
39 nowadays within easy reach of non-statisticians thanks to their implementation in user-
40 friendly software packages. Here, we would like to stress once more that, for all of the
41 large variety of available procedures and their complications, the basic logic of hypothesis
42 testing and estimation is essentially the same that we illustrated with very simple exam-
43 ples of frequency count data in the first sections of this article. It is not essential to know
44 mathematical details of all the techniques in order to apply them, but it is important
45 to understand the basic principles of hypothesis testing and estimation; the assumptions
46 of a test, its null hypothesis, and the meaning of a p-value; and to make sure that the
47 assumptions are met by the data and that the research question can be translated into a
48 meaningful null hypothesis. And, of course, the linguistic interpretation of the statistical
49 results is at least as crucial as the correctness of the methods applied.

1 We have focused here on statistical inference for hypothesis testing and estimation, as
2 applied to corpus data. This is only a part, albeit a fundamental one, of the role that
3 statistical methods play in corpus-related disciplines today. For a survey of statistical
4 procedures used for *exploratory* purposes (i.e., as an aid in uncovering interesting patterns
5 in the data), see articles 40 and 42. Statistical methods also play a very important role
6 as *modeling* tools for machine learning techniques applied to natural language (article 41)
7 and more generally in so-called empirical natural language processing (see, e.g., article
8 50 on machine translation, and Manning/Schütze 1999 for an introduction to statistical
9 NLP).

10 9 Directions for further study

11 A much more in-depth introduction to the statistical inference methods appropriate for
12 count data that we discussed here is provided by Agresti (1996) or, at a more technical
13 level, Agresti (2002). There is, of course, a vast number of introductory general statistics
14 books. DeGroot/Schervish (2002) present a particularly thorough and clear introduc-
15 tion, although it requires at least a basic mathematical background. Among the less
16 technical introductions, we recommend the one by Hays (1994), a book that provides
17 non-mathematical but rigorous explanations of the most important notions of statistical
18 inference (although it focuses on the statistical methods for the analysis of experimental
19 results, which are only partially relevant to corpus work). There is also a wealth of statis-
20 tics “cookbooks” that illustrate when to apply a certain technique, how to apply it, and
21 how to interpret the results. These are often and usefully linked to a specific statistical
22 software package. For example, Dalgaard (2002) is an introduction to running various
23 standard statistical procedures in R (see below).

24 There are a few older introductions to statistics explicitly geared towards linguists.
25 The one by Woods/Fletcher/Hughes (1986) is a classic, whereas the one by Butler (1985)
26 has the advantage of being now freely available on the Web:

27 <http://www.uwe.ac.uk/hlss/llas/statistics-in-linguistics/bkindex.shtml>

28 Older introductions tend to focus on techniques that are more relevant to psycholin-
29 guistics, phonetics and language testing than to corpus analysis. Oakes (1998) presents a
30 survey of applications of statistics in corpus studies that trades depth for wider breadth
31 of surveyed applications and methods. It is likely that, with the growing interest in cor-
32 pora and statistical approaches to linguistics in general, the next few years will see the
33 appearance of more statistics textbooks targeting corpus linguists.

34 There is nowadays a large number of statistical software packages to choose from. We
35 recommend R:

36 <http://www.r-project.org/>

37 R supports an impressive range of statistical procedures and, being open-source and
38 available free of charge, it is attracting a growing community of developers who add new
39 functionalities, including some that are of interest to corpus linguists. These extensions
40 range from advanced data visualization techniques to modules explicitly targeting corpus
41 work, such as the *corpora* library developed by the authors of this article:

42 <http://purl.org/stefan.evert/SIGIL>

1 The corpora library (also free and open-source) provides support for carrying out the
2 statistical analyses described in this article (the Web site has a tutorial that shows how to
3 run them), as well as several sample data sets. There is an increasing number of introduc-
4 tory textbooks with concrete R examples, and we know of several R-based books focusing
5 on statistical methods in linguistics that are currently in preparation. Shravan Vasishth
6 has written (and is constantly updating) an online book aimed at (psycho-)linguists that
7 introduces statistics in R through a simulation approach. This book is freely available
8 (under a Creative Commons license) from:

9 <http://www.ling.uni-potsdam.de/~vasishth/SFLS.html>

10 Finally, Wulff (2005) provides a survey of online statistics facilities.

11 References

- 12 Agresti, Alan (1996), *An Introduction to Categorical Data Analysis*. Chichester: Wiley.
- 13 Agresti, Alan (2002), *Categorical Data Analysis, second edition*. Chichester: Wiley.
- 14 Baayen, Harald (2001), *Word Frequency Distributions*. Dordrecht: Kluwer.
- 15 Biber, Douglas/Conrad, Susan/Reppen, Randi (1998), *Corpus Linguistics*. Cambridge:
16 Cambridge University Press.
- 17 Biber, Douglas/Johansson, Stig/Leech, Geoffrey/Conrad, Susan/Finegan, Edward (1999),
18 *Longman Grammar of Spoken and Written English*. Harlow, UK: Pearson Education.
- 19 Butler, Christopher (1985), *Statistics in Linguistics*. Oxford: Blackwell.
- 20 Chomsky, Noam (1986), *Knowledge of Language: Its Nature, Origins, and Use*. New York:
21 Praeger.
- 22 Church, Kenneth (2000), Empirical estimates of adaptation: the chance of two Noriegas
23 is closer to $p/2$ than p^2 . In: *Proceedings of the 17th Conference on Computational*
24 *Linguistics*, 180-186.
- 25 Culicover, Peter/Jackendoff, Ray (2005), *Simpler Syntax*. Oxford: Oxford University
26 Press.
- 27 Dalgaard, Peter (2002), *Introductory Statistics with R*. New York: Springer.
- 28 DeGroot, Morris/Schervish, Mark (2002), *Probability and Statistics, third edition*. Boston:
29 Addison-Wesley.
- 30 Evert, Stefan (2006). How random is a corpus? The library metaphor. In: *Zeitschrift für*
31 *Anglistik und Amerikanistik* 54(2), 177-190.
- 32 Hays, William (1994), *Statistics, fifth edition*. New York: Harcourt Brace.
- 33 Gries, Stefan Th. (2005), Null-hypothesis significance testing of word frequencies: A follow-
34 up on Kilgarriff. In: *Corpus Linguistics and Linguistic Theory* 1, 277-294.
- 35 Manning, Christopher/Schütze, Hinrich (1999), *Foundations of Statistical Natural Lan-*
36 *guage Processing*. Cambridge (Mass.): MIT Press.

- 1 McEnery, Tony/Wilson, Andrew (2001), *Corpus Linguistics, second edition*. Edinburgh:
2 Edinburgh University Press.
- 3 Oakes, Michael (1998), *Statistics for Corpus Linguistics*. Edinburgh: Edinburgh University
4 Press.
- 5 Pearson, Egon (1990), *'Student': A Statistical Biography of William Sealy Gosset*. Oxford:
6 Clarendon Press.
- 7 Rainer, Franz (2003), Studying restrictions on patterns of word-formation by means of the
8 Internet. In: *Rivista di Linguistica* 15, 131-139.
- 9 Schütze, Carson (1996), *The Empirical Base of linguistics: Grammaticality Judgments
10 and Linguistic Methodology*. Chicago: University of Chicago Press.
- 11 Woods, Anthony/Fletcher, Paul/Hughes, Arthur (1986), *Statistics in Language Studies*.
12 Cambridge: CUP.
- 13 Wulff, Stefanie (2005), Online statistics labs. In: *Corpus Linguistics and Linguistic Theory*
14 1, 303-308.