

Shravan Vasishth, Daniel Schad, Audrey Bürki, Reinhold Kliegl

Linear Mixed Models in Linguistics and Psychology: A Comprehensive Introduction

Dedicated to ...

Contents

Preface	vii
0.1 Prerequisites	viii
0.2 How to read this book	viii
0.3 Online materials	ix
0.4 Software needed	ix
0.5 Acknowledgements	x
 About the Authors	 xi
 I Foundational ideas	 1
 1 Some important facts about distributions	 3
1.1 Discrete random variables: An example using the Binomial distribution	3
1.1.1 The mean and variance of the Binomial dis- tribution	5
1.1.2 What information does a probability distri- bution provide?	7
1.2 Continuous random variables: An example using the Normal distribution	11
1.3 Bivariate/multivariate distributions	15
1.4 Summary of useful R functions relating to univari- ate distributions	19
1.5 *Summary of random variable theory	19
1.6 Further reading	21
1.7 Exercises	22
1.7.1 Practice using the <code>pnorm</code> function	22
1.7.2 Practice using the <code>qnorm</code> function	22
1.7.3 Practice using <code>qt</code>	23
1.7.4 Maximum likelihood estimation 1	23

1.7.5	Maximum likelihood estimation 2	24
1.7.6	Generating bivariate data	25
1.7.7	Generating multivariate data	25
2	Hypothetical repeated sampling and the t-test	27
2.1	The central limit theorem using simulation	27
2.2	Three examples of the sampling distribution . . .	32
2.3	The confidence interval, and what it's good for .	34
2.4	Hypothesis testing: The one sample t-test	36
2.4.1	The one-sample t-test	37
2.4.2	Type I, II error, and power	45
2.4.3	How to compute power for the one-sample t-test	49
2.4.4	The p-value	52
2.4.5	Type M and S error in the face of low power	56
2.4.6	Searching for significance	60
2.5	The two-sample t-test vs. the paired t-test	63
2.5.1	Common mistakes involving the t-test . .	67
2.6	Exercises	74
2.6.1	Computing the p-value	74
2.6.2	Computing the t-value	74
2.6.3	Type I and II error	74
3	Linear models and linear mixed models	75
3.1	From the t-test to the linear (mixed) model . . .	75
3.2	Sum coding	82
3.3	Checking model assumptions	84
3.4	From the paired t-test to the linear mixed model	87
3.5	Linear mixed models	94
3.5.1	Model type 1: Varying intercepts	99
3.5.2	The formal statement of the varying inter- cepts model	100
3.5.3	Model type 2: Varying intercepts and slopes, without a correlation	102
3.5.4	Model type 3: Varying intercepts and vary- ing slopes, with correlation	108
3.6	Shrinkage in linear mixed models	111

3.7	Summary	114
3.8	Exercises	114
3.8.1	By-subjects t-test	115
3.8.2	Fitting a linear mixed model	115
3.8.3	t-test vs. linear mixed model	115
3.8.4	Power calculation using power.t.test	116
3.8.5	Residuals	116
3.8.6	Understanding contrast coding	117
3.8.7	Understanding the fixed-effects output	117
3.8.8	Understanding the null hypothesis test	117
4	Hypothesis testing using the likelihood ratio test	119
4.1	The likelihood ratio test: The theory	119
4.2	A practical example using simulated data	124
4.3	A real-life example: The English relative clause data	126
4.4	Exercises	129
4.4.1	Chinese relative clauses	129
4.4.2	Agreement attraction in comprehension	130
4.4.3	The grammaticality illusion	132
5	Using simulation to understand your model	135
5.1	A reminder: The maximal linear mixed model	135
5.2	Obtain estimates from a previous study	136
5.3	Decide on a range of plausible values of the effect size	138
5.4	Extract parameter estimates	140
5.5	Define a function for generating data	142
5.5.1	Generate a Latin-square design	142
5.5.2	Generate data row-by-row	143
5.6	Repeated generation of data to compute power	150
5.7	What you can now do	153
5.8	Exercises	155
5.8.1	Drawing a power curve given a range of effect sizes	155
5.8.2	Power and log-transformation	155

5.8.3	Evaluating models by generating simulated data	155
5.8.4	Using simulation to check parameter recovery	156
5.8.5	Sample size calculations using simulation .	156
6	Important distributions	157

Preface

This book is intended to be a relatively complete introduction to the application of linear mixed models in areas related to linguistics and psychology; throughout, we use the programming language R. Our target audience is cognitive scientists (e.g., linguists and psychologists) who carry out behavioral experiments, and who are interested in learning the foundational ideas behind modern statistical methodology from the ground up and in a principled manner.

Many excellent introductory textbooks already exist that discuss data analysis in great detail. Our book is different from existing books in two respects. First, our main focus is on showing how to analyze data from planned experiments involving repeated measures; this type of experimental data involves complexities that are distinct from the problems one encounters when analyzing observational data. We provide many examples of data-sets involving eyetracking (visual world and reading), self-paced reading, event-related potentials, reaction time, acceptability rating judgements, speeded grammaticality judgements, and question-response accuracies. Second, from the very outset, we stress a particular workflow that has as its centerpiece simulating data; we aim to teach a philosophy that involves thinking hard about the assumed underlying generative process, **even before the data are collected**. The data analysis approach that we hope to teach through this book involves a cycle of experiment design analysis and model validation using simulated data.

0.1 Prerequisites

This book assumes high school arithmetic and algebra. We also expect that the reader already knows basic constructs in the programming language R (R Core Team, 2019), such as writing for-loops. For newcomers to R, we provide a quick introduction in the appendix that covers all the constructs used in the book. For those lacking background in R, there are many good online resources on R that they can consult as needed. Examples are: R for data science¹, and Efficient R programming².



provide comprehensive book recommendations

0.2 How to read this book

The chapters in this book are intended to be read in sequence, but during the first pass through the book, the reader should feel free to skip the sections marked with an asterisk. These sections provide a more formal development that will be useful when the reader transitions to more advanced textbooks like Gelman et al. (2014).

to-do: add a Mackay type chapter ordering for different scenarios.

¹<https://r4ds.had.co.nz/>

²<https://csgillespie.github.io/efficientR/>

0.3 Online materials

The entire book, including all data and source code, is available online from https://github.com/vasishth/Freq_CogSci³. Solutions to the exercises will be provided (to-do).

0.4 Software needed

Before you start, please install

- R⁴ (and RStudio⁵ or any other IDE)
- The R packages MASS, dplyr, purrr, readr, extraDistr, ggplot2:
 - They can be installed in the usual way:
`install.packages(c("MASS", "dplyr", "purrr",
"readr", "extraDistr", "ggplot2"))`.

In every R session, we'll need to set a seed (this ensures that the random numbers are always the same).

```
set.seed(42)
library(MASS)
## be careful to load dplyr after MASS
library(dplyr)
library(purrr)
library(readr)
library(extraDistr)
library(ggplot2)
```

³https://github.com/vasishth/Freq_CogSci

⁴<https://cran.r-project.org/>

⁵<https://www.rstudio.com/>

0.5 Acknowledgements

We are grateful to the many generations of students at the University of Potsdam, various summer schools at ESSLLI, the LOT winter school, other short courses we have taught at various institutions, and the annual summer school on Statistical Methods for Linguistics and Psychology (SMLP) at the University of Potsdam. The participants in these courses helped us considerably in improving the material presented here. We are also grateful to members of Vasishth lab for comments on earlier drafts of this book.

Shravan Vasishth, Daniel Schad, Audrey Bürki, Reinhold Kliegl,
Potsdam, Germany

About the Authors

Shravan Vasishth (<http://vasishth.github.io>) is professor of Psycholinguistics at the University of Potsdam, Germany. He holds the chair for Psycholinguistics and Neurolinguistics (Language Processing). After completing his Bachelors degree in Japanese from Jawaharlal Nehru University, New Delhi, India, he spent five years in Osaka, Japan, studying Japanese and then working as a translator in a patent law firm in Osaka. He completed an MS in Computer and Information Science (2000-2002) and a PhD in Linguistics (1997-2002) from the Ohio State University, Columbus, USA, and an MSc in Statistics (2011-2015) from the School of Mathematics and Statistics, University of Sheffield, UK. He is a professional member of the Royal Statistical Society (GradStat ID: 128307), a member of the International Society for Bayesian Analysis, and a lifetime member of the Linguistic Society of America (LSA). He is on the editorial board of the Linguistic Society of America flagship journal *Language* as their statistics consultant for journal submissions. His research focuses on computational modeling of sentence processing in unimpaired and impaired populations, and the application of mathematical, computational, experimental, and statistical methods (particularly Bayesian methods) in linguistics and psychology. He runs an annual summer school, Statistical Methods in Linguistics and Psychology (SMLP): vasishth.github.io/smlp. He regularly teaches short courses on statistical data analysis (Bayesian and frequentist methods).

Daniel J. Schad (<https://danielschad.github.io/>) is an assistant professor in the department of Cognitive Science and AI at Tilburg University. He studied Psychology at the University of Potsdam, Germany, and at the University of Michigan, Ann Arbor, USA. He did a PhD in Cognitive Psychology at the University of

Potsdam, working on computational models of eye-movement control and on mindless reading. He then did a five-year post-doc in the novel field of Computational Psychiatry at the Charité, Universität Berlin, Germany (partly also at the University of Potsdam), with research visits at the ETH Zürich, Switzerland, and the University College London, UK, working on model-free and model-based decision-making and Pavlovian-instrumental transfer in alcohol dependence, and on the cognitive and brain mechanisms underlying Pavlovian conditioning. He has worked as a postdoctoral researcher at the University of Potsdam, conducting research on quantitative methods in Cognitive Science, including contrasts, properties of significance tests, Bayesian Workflow, and Bayes factor analyses.

Audrey Bürki (<https://audreyburki.github.io/Website/>) leads a research group at the University of Potsdam, Germany. Her research interests are... to-do

Reinhold Kliegl (<https://www.uni-potsdam.de/de/trainingwissenschaft/mitarbeiter/rkliegl.html>) is a senior professor in the Division of Training and Movement Science, at the University of Potsdam, Germany. His research interests are... to-do



Part I

Foundational ideas



1

Some important facts about distributions

In linguistics and psychology, typical data-sets involve either *discrete* dependent measures such as acceptability ratings on a Likert scale (for example, ranging from 1 to 7), and binary grammaticality judgements, or *continuous* dependent measures such as reading times or reaction times in milliseconds and EEG signals in microvolts. Whenever we fit a model using one of these types of dependent measures, we make some assumptions about how these measurements were generated. In particular, we usually assume that our observed measurements are coming from a particular *probability mass function* (when the data are discrete in nature) *probability density function* (when the data are continuous). Behind these PMFs and PDFs lies the concept of a random variable; as will become apparent in this chapter, it is extremely useful to be able to think about data in terms of the random variable assumed. We consider the two cases, discrete and continuous, separately.

1.1 Discrete random variables: An example using the Binomial distribution

Imagine that our data come from a grammaticality judgement task, and that the responses from participants are a sequence of 1's and 0's, where 1 represents the judgment “grammatical”, and 0 represents the judgement “ungrammatical”. Assume also that each response, coded as 1 or 0, is generated independently from the others. We can simulate such a sequence of 1s and 0s in R. Here is a case where we run the same experiment 20 times (the sample size is 10 each time).

```
## [1] 7 7 4 7 6 5 6 3 6 6 5 6 7 4 5 7 8 3 5 5
```

The number of successes in each of the 20 simulated experiments above is being generated by a *discrete random variable* Y with a probability distribution $p(Y)$ called the **Binomial distribution**.¹

For discrete random variable, the probability distribution $p(Y)$ is called a **probability mass function** (PMF). The PMF defines the probability of each possible outcome. In the above example, with $n = 10$ trials, there are 11 possible outcomes: $0, \dots, 10$ successes. Which of these outcomes is most probable depends on a parameter in the Binomial distribution that represents the probability of success. We will call this parameter θ . The left-hand side plot in Figure 1.1 shows an example of a Binomial PMF with 10 trials and the parameter θ with value 0.5. Setting θ to 0.5 leads to a PMF where the most probable outcome is 5 successes out of 10. If we had set θ to, say 0.1, then the most probable outcome would be 1 success out of 10; and if we had set θ to 0.9, then the most probable outcome would be 9 successes out of 10.

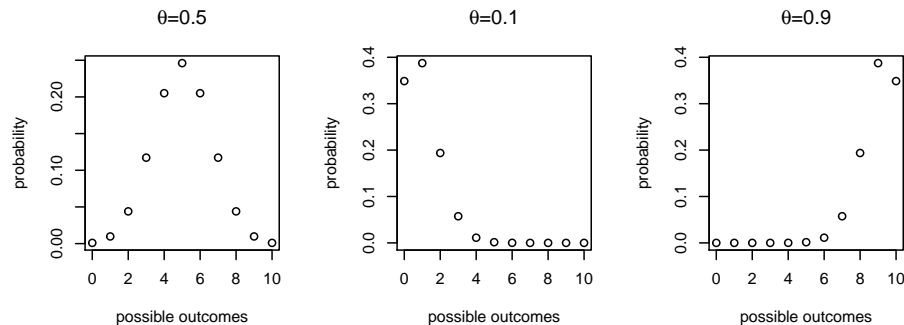


FIGURE 1.1: Probability mass functions of a binomial distribution assuming 10 trials, with 50%, 10%, and 90% probability of success.



to-do bar or line graphs above, instead of points

¹When an experiment consists of only a single trial (i.e., we can have a total number of only 0 or 1 successes), $p(Y)$ is called a **Bernoulli distribution**.

The probability mass function for the binomial is written as follows.

$$\text{Binomial}(k|n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (1.1)$$

Here, n represents the total number of trials, k the number of successes, and θ the probability of success. The term $\binom{n}{k}$, pronounced n -choose- k , represents the number of ways in which one can choose k successes out of n trials. For example, 1 success out of 10 can occur in 10 possible ways: the very first trial could be a 1, the second trial could be a 1, etc. The term $\binom{n}{k}$ expands to $\frac{n!}{k!(n-k)!}$. In R, it is computed using the function `choose(n,k)`, with n and k representing positive integer values.

1.1.1 The mean and variance of the Binomial distribution

It is possible to analytically compute the mean and variance of the PMF associated with the Binomial random variable Y . Without getting into the details of how these are derived mathematically, we just state here that the mean of Y (also called the expectation, conventionally written $E[Y]$) and variance of Y (written $\text{Var}(Y)$) of a Binomial distribution with parameter θ and n trials are $E[Y] = n\theta$ and $\text{Var}(Y) = n\theta(1 - \theta)$, respectively.

Of course, we always know n (because we decide on the number of trials ourselves), but in real experimental situations we never know the true value of θ . But θ can be estimated from the data. From the observed data, we can compute the estimate of θ , $\hat{\theta} = k/n$. The quantity $\hat{\theta}$ is the observed proportion of successes, and is called the **maximum likelihood estimate** of the true (but unknown mean). Once we have estimated θ in this way, we can also obtain an estimate (also a maximum likelihood estimate) of the variance by computing $n\hat{\theta}(1 - \hat{\theta})$. These estimates are then used for statistical inference.

What does the term “maximum likelihood estimate” mean? The term **likelihood** refers to the value of the Binomial distribution function for a particular value of θ , once we have observed some

data. For example, suppose you record $n = 10$ trials, and observe $k = 7$ successes. What is the probability of observing 7 successes out of 10? We need the binomial distribution to compute this value:

$$\text{Binomial}(k = 7 | n = 10, \theta) = \binom{10}{7} \theta^7 (1 - \theta)^{10-7} \quad (1.2)$$

Once we have observed the data, both n and k are fixed. The only variable in the above equation now is θ : the above function is now only dependent on the value of θ . When the data are fixed, the probability mass function is only dependent on the value of the parameter θ , and is called a **likelihood function**. It is therefore often expressed as a function of θ :

$$p(y|\theta) = p(k = 7, n = 10|\theta) = \mathcal{L}(\theta)$$

The vertical bar notation above should be read as saying that, given some data y (which in the binomial case will be k “successes” in n trials), the function returns a value for different values of θ .

If we now plot this function for all possible values of θ , we get the plot shown in Figure 1.2.



DS comment: do we want to show the code for computing all likelihood values? (maybe this comes later?)

What is important about this plot is that it shows that, given the data, the maximum point is at the point 0.7, which corresponds to the estimated mean using the formula shown above: $k/n = 7/10$. Thus, the maximum likelihood estimate (MLE) gives us the most likely value that the parameter θ given the data. It is crucial to note here that the phrase “most likely” here does not mean that the MLE from a *particular* sample of data invariably gives us an accurate estimate of θ . For example, if we run our experiment for 10 trials and get 1 success out of 10, the MLE is 0.10. We could have happened to observe only one success out of ten even if the true θ were 0.5. The MLE would however give an accurate estimate of the true parameter as n approaches infinity.

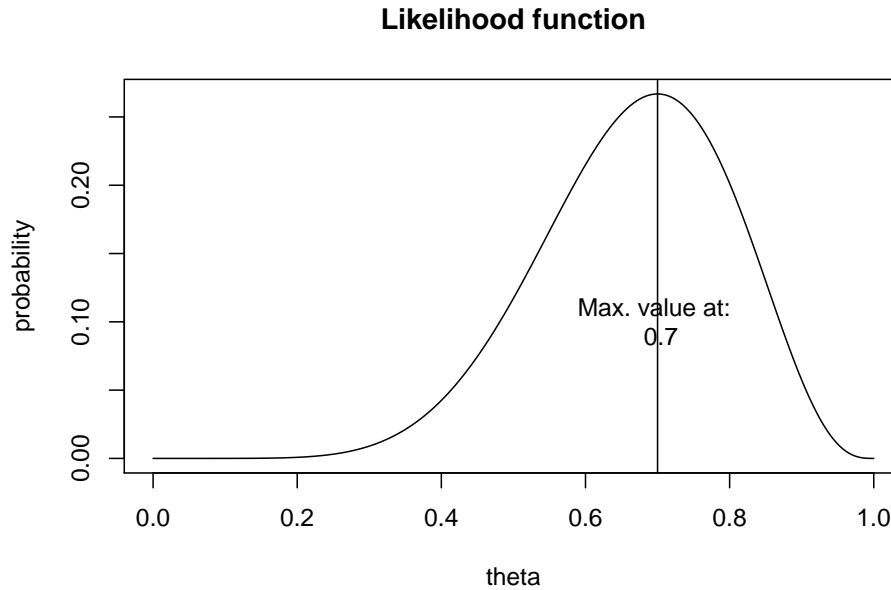


FIGURE 1.2: The likelihood function for 7 successes out of 10.

1.1.2 What information does a probability distribution provide?

What good is a probability mass function? We consider this question next.

1.1.2.1 Compute the probability of a particular outcome (discrete case only)

The Binomial distribution shown in Figure 1.1 already shows the probability of each possible outcome under a different value for θ . In R, there is a built-in function that allows us to calculate the probability of k successes out of n , given a particular value of k (this number constitutes our data), the number of trials n , and given a particular value of θ ; this is the `dbinom` function. For example, the probability of 5 successes out of 10 when θ is 0.5 is:

```
dbinom(5, size = 10, prob = 0.5)
```

```
## [1] 0.2461
```

The probabilities of success when θ is 0.1 or 0.9 can be computed by replacing 0.5 above by each of these probabilities. One can just do this by giving `dbinom` a vector of probabilities:

```
dbinom(5, size = 10, prob = c(0.1, 0.9))
```

```
## [1] 0.001488 0.001488
```

1.1.2.2 Compute the cumulative probability of k or less (more) than k successes

Using the `dbinom` function, we can compute the cumulative probability of obtaining 1 or less, 2 or less successes etc. This is done through a simple summation procedure:

```
## the cumulative probability of obtaining 0, 1, or
## 2 successes out of 10, with theta=0.5:
dbinom(0, size = 10, prob = 0.5) + dbinom(1, size = 10,
  prob = 0.5) + dbinom(2, size = 10, prob = 0.5)
```

```
## [1] 0.05469
```

Mathematically, we could write the above summation as:

$$\sum_{k=0}^2 \binom{n}{k} \theta^k (1-\theta)^{n-k} \quad (1.3)$$

An alternative to the cumbersome addition in the R code above is this more compact statement, which closely mimics the above mathematical expression:

```
sum(dbinom(0:2, size = 10, prob = 0.5))
```

```
## [1] 0.05469
```

R has a built-in function called `pbinom` that does this summation for us. If we want to know the probability of 2 or less successes as in the above example, we can write:

```
pbinom(2, size = 10, prob = 0.5, lower.tail = TRUE)
```

```
## [1] 0.05469
```

The specification `lower.tail=TRUE` ensures that the summation goes from 2 to numbers smaller than 2 (which lie in the lower tail of the distribution in Figure 1.1). If we wanted to know what the probability is of obtaining 2 or more successes out of 10, we can set `lower.tail` to `FALSE`:

```
pbinom(2, size = 10, prob = 0.5, lower.tail = FALSE)
```

```
## [1] 0.9453
```

The cumulative distribution function or CDF can be plotted by computing the cumulative probabilities for any value k or less than k , where k ranges from 0 to 10 in our running example. The CDF is shown in Figure 1.3.

1.1.2.3 Compute the inverse of the cumulative distribution function (the quantile function)

We can also find out the value of the variable k (the quantile) such that the probability of obtaining k or less than k successes is some specific probability value p . If we switch the x and y axes of Figure 1.3, we obtain another very useful function, the inverse CDF.

The inverse of the CDF (known as the quantile function in R because it returns the quantile, the value k) is available in R as the function `qbinom`. The usage is as follows: to find out what the value k of the outcome is such that the probability of obtaining k or less successes is 0.37, type:

```
qbinom(0.37, size = 10, prob = 0.5)
```

```
## [1] 4
```

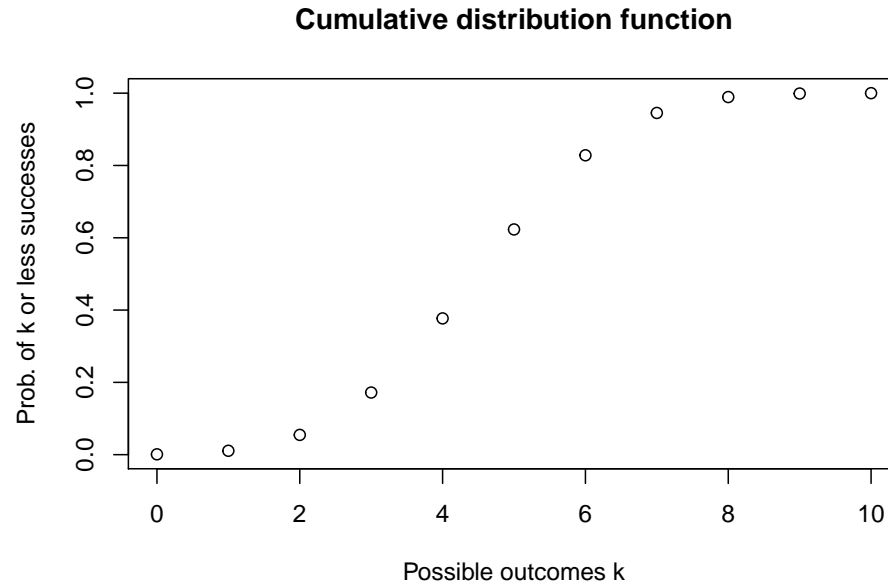


FIGURE 1.3: The cumulative distribution function for a binomial distribution assuming 10 trials, with 50% probability of success.

1.1.2.4 Generate random data from a $\text{Binomial}(n, \theta)$ distribution

We can generate random simulated data from a Binomial distribution by specifying the number of trials and the probability of success θ . In R, we do this as follows:

```
rbinom(10, size = 1, prob = 0.5)
```

```
## [1] 1 0 1 1 0 1 0 1 0 1
```

The above code generates a sequences of 1's and 0's. Repeatedly run the above code; you will get different sequences each time. For each generated sequence, one can calculate the number of successes by just summing up the vector, or computing its mean and multiplying by the number of trials, here 10:

```
y <- rbinom(10, size = 1, prob = 0.5)
mean(y) * 10
```

```
## [1] 6
```

```
sum(y)
```

```
## [1] 6
```

1.2 Continuous random variables: An example using the Normal distribution

We will now revisit the idea of a random variable using a continuous distribution. Imagine that you have a vector of reading time data y measured in milliseconds and coming from a Normal distribution. The probability density function (PDF) of the Normal distribution is defined as follows:

$$Normal(y|\mu, \sigma) = f(y) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y-\mu)^2}{2\sigma^2}\right) \quad (1.4)$$

Here, μ is the mean, and σ is the standard deviation of the Normal distribution that the reading times have been sampled from.

We can visualize the Normal distribution for particular values of μ and σ as a PDF (using `dnorm`), a CDF (using `pnorm`), and the inverse CDF (using `qnorm`). See Figure 1.4.

In the figure, the PDF gives us the density for each possible value of our data y ; the CDF tells us the probability of observing a value like y or some value less than that (written: $P(Y < y)$); and the inverse CDF gives us the quantile y such that $P(Y < y)$ is some specific value between 0 and 1. These are three different ways of looking at the same information.

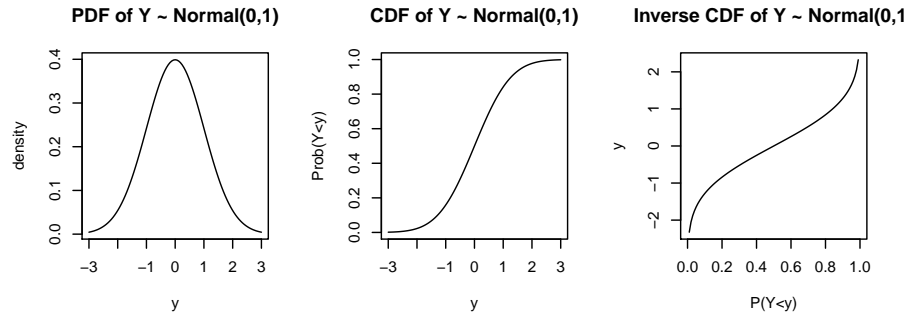


FIGURE 1.4: The PDF, CDF, and inverse CDF for the $Normal(\mu = 0, \sigma = 1)$.

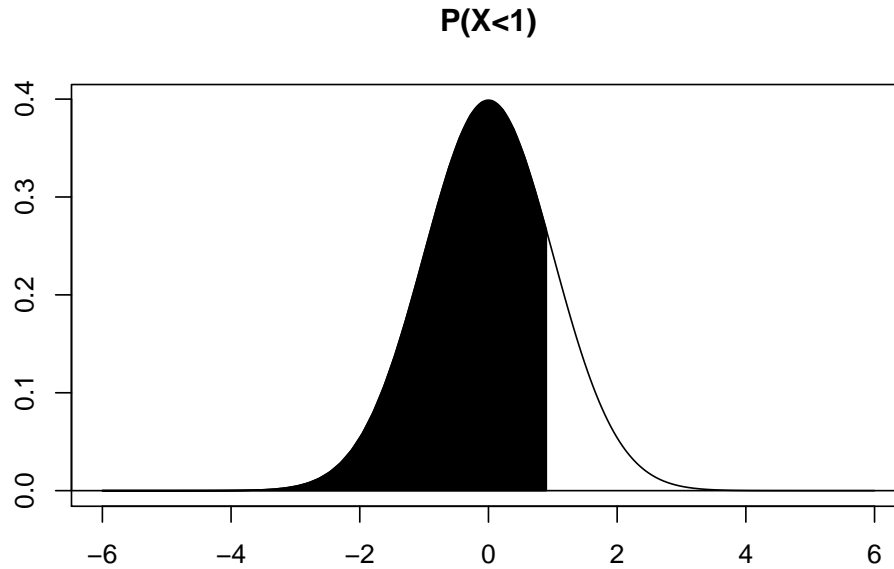
One important fact about the normal distribution is that 95% of the probability mass is covered by approximately plus/minus 1.96 times the standard deviation about the mean. Thus, the range $\mu \pm 1.96 \times \sigma$ will cover approximately 95% of the area under the curve. We will approximate this by talking about $\mu \pm 2 \times \sigma$.

As in the discrete example, the PDF, CDF, and inverse of the CDF allow us to ask questions like:

- **What is the probability of observing values between a and b from a Normal distribution with mean μ and standard deviation σ ?** We can compute the probability of the random variable lying between 1 and minus infinity:

```
pnorm(1, mean = 0, sd = 1) - pnorm(-Inf, mean = 0,
sd = 1)
```

```
## [1] 0.8413
```

Notice here that the probability of any point value in a PDF is always 0. This is because the probability in a continuous probability distribution is the area under the curve, and the area at any point on the x-axis is always 0. The implication here is that we can only ask about probabilities between two different points; e.g., the probability that Y lies between a and b , or $P(a < Y < b)$. Also, notice that $P(a < Y < b)$ and $P(a \leq Y \leq b)$ will be the same probability, because of the fact that $P(Y = a)$ or $P(Y = b)$ both equal 0.

- **What is the quantile q such that the probability is p of observing that value q or something less (or more) than it?** For example, we can work out the quantile q such that the probability of observing q or something less than it is 0.975, in the $\text{Normal}(500,100)$ distribution. Formally, we would write this as $P(Y < a)$.

```
qnorm(0.975, mean = 500, sd = 100)
```

```
## [1] 696
```

The above output says that the probability that the random variable is less than $q = 695$ is 97.5%.

- **Generating simulated data.** Given a vector of n independent and identically distributed data y , i.e., given that each data point is being generated independently from $Y \sim \text{Normal}(\mu, \sigma)$ for some values of the parameters, the maximum likelihood estimates for the expectation and variance are

$$\bar{y} = \frac{\sum_{i=1}^n y_i}{n} \quad (1.5)$$

$$\text{Var}(y) = \frac{\sum_{i=1}^n (y_i - \bar{y})^2}{n} \quad (1.6)$$

For example, you could generate 10 data points using the `rnorm` function, and then compute the mean and variance from the simulated data:

```
y <- rnorm(10, mean = 500, sd = 100)
mean(y)
```

```
## [1] 482.2
```

```
var(y)
```

```
## [1] 13365
```

Again, depending on the sample size, the sample mean and sample variance may or may not be close to the true values of the respective parameters, despite the fact that these are maximum likelihood estimates.

One other important detail about probability distributions is that one can have a joint distribution defined for more than one random variable. We consider this situation next, considering the continuous case only.

1.3 Bivariate/multivariate distributions

So far, we have only discussed univariate distributions. It is also possible to specify distributions with two or more dimensions. For example, consider the bivariate case: two normal random variables U_0 and U_1 , each of which coming from, for example, a Normal(0,1) distribution, with some correlation ρ between the random variables.

A bivariate distribution for two random variables U_0 and U_1 , each of which comes from a normal distribution, is expressed in terms of the means and standard deviations of each of the two distributions, and the correlation ρ between them. The standard deviations and correlation are expressed in a special form of a 2×2 matrix called a variance covariance matrix Σ . If ρ_u is the correlation between the two random variables, and σ_{u0} and σ_{u1} the respective standard deviations, the variance covariance matrix is written as:

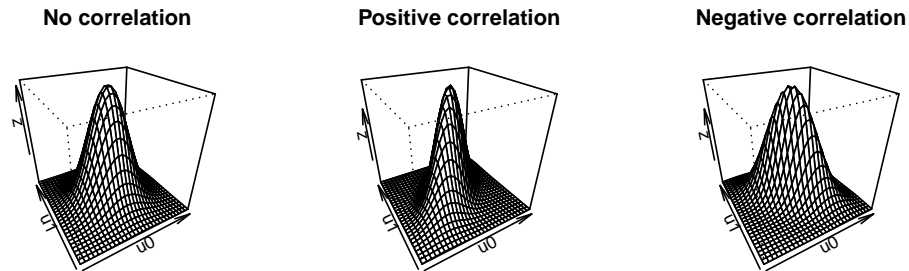
$$\Sigma = \begin{pmatrix} \sigma_{u0}^2 & \rho_u \sigma_{u0} \sigma_{u1} \\ \rho_u \sigma_{u0} \sigma_{u1} & \sigma_{u1}^2 \end{pmatrix} \quad (1.7)$$

The off-diagonals of this matrix contain the covariance between U_0 and U_1 .

The joint distribution of U_0 and U_1 is defined as follows:

$$\begin{pmatrix} U_0 \\ U_1 \end{pmatrix} \sim \mathcal{N}_2 \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma \right) \quad (1.8)$$

A visualization will help. The figure below shows a bivariate distribution with correlation zero, a positive and a negative correlation.



In this book, we will make use of such multivariate distributions a lot, and it will soon become important to know how to generate simulated bivariate data that is correlated. Suppose we want to generate 100 correlated pairs of data, with correlation $\rho = 0.6$. The two random variables have mean 0, and standard deviations 5 and 10 respectively.

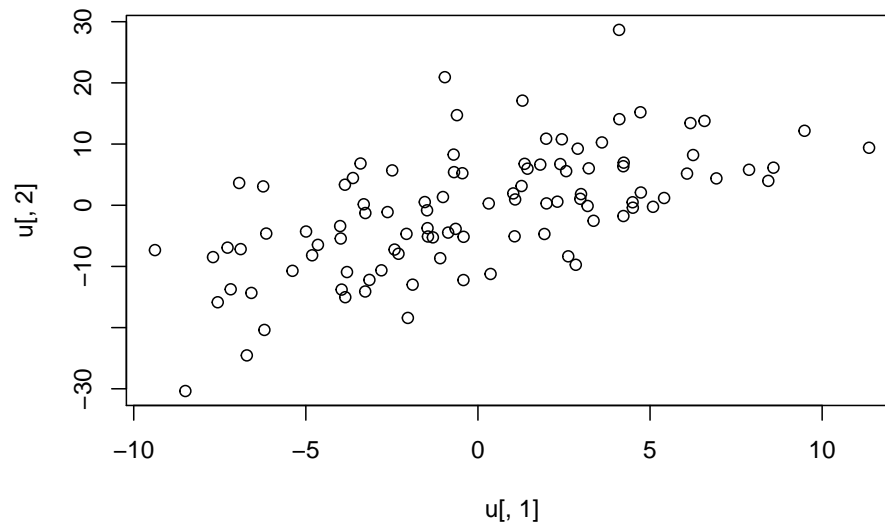
Here is how we would generate such data. First, define a variance covariance matrix; then, use the multivariate analog of the `rnorm` function, `mvrnorm`, to generate 100 data-points.

```
library(MASS)
## define a variance-covariance matrix:
Sigma <- matrix(c(5^2, 5 * 10 * 0.6, 5 * 10 * 0.6,
  10^2), byrow = FALSE, ncol = 2)
## generate data:
u <- mvrnorm(n = 100, mu = c(0, 0), Sigma = Sigma)
head(u)
```

```
##      [,1]      [,2]
## [1,] -0.453    5.233
## [2,]  4.222    6.372
## [3,]  3.602   10.269
## [4,] -6.151   -4.627
## [5,] -3.412    6.823
## [6,] -2.036  -18.415
```

A plot confirms that the data are positively correlated.

```
plot(u[, 1], u[, 2])
```



As an exercise, try changing the correlation to 0 or to -0.6 , and then plot the bivariate distribution that results.

One final useful thing to notice about the variance-covariance matrix is that it can be decomposed into the component standard deviations and an underlying correlation matrix. For example, consider the matrix above:

```
Sigma
```

```
##      [,1] [,2]
## [1,]  25  30
## [2,]  30 100
```

One can decompose the matrix as follows. The matrix can be seen as the product of a diagonal matrix of the standard deviations and the correlation matrix:

```
## sds:
(sds <- c(5, 10))
```

```
## [1] 5 10
```

```
## diagonal matrix:
(sd_diag <- diag(sds))
```

```
##      [,1] [,2]
## [1,] 5    0
## [2,] 0    10
```

```
## correlation matrix:
(corrmatrix <- matrix(c(1, 0.6, 0.6, 1), ncol = 2))
```

```
##      [,1] [,2]
## [1,] 1.0 0.6
## [2,] 0.6 1.0
```

Given these two matrices, one can reassemble the variance covariance matrix:

```
sd_diag %*% corrmatrix %*% sd_diag
```

```
##      [,1] [,2]
## [1,] 25  30
## [2,] 30 100
```

There is a built-in convenience function, `sdcor3cov` in the `SIN` package that does this calculation, taking the vector of standard deviations (not the diagonal matrix) and the correlation matrix to yield the variance covariance matrix:

```
SIN::sdcor2cov(stddev = sds, corr = corrmatrix)
```

```
##      [,1] [,2]
## [1,] 25  30
## [2,] 30 100
```

We will be using this function a lot when simulating data from linear mixed models.

1.4 Summary of useful R functions relating to univariate distributions

Table 1.1 summarizes the different functions relating to univariate PMFs and PDFs, using the Binomial and Normal as examples.

TABLE 1.1: Important R functions relating to two univariate random variables, the Normal and the Binomial.

	Discrete	Continuous
Example:	Binomial(n, θ)	Normal(μ, σ)
Likelihood function	dbinom	dnorm
Prob $Y=y$	dbinom	always 0
Prob $Y \geq y, Y \leq y, y_1 < Y < y_2$	pbinom	pnorm
Inverse CDF	qbinom	qnorm
Generate simulated data	rbinom	rnorm

Other distributions, such as the t-distribution, the Uniform, Exponential, Gamma, etc. Beta, etc., have their own set of d-p-q-r functions in R. The appendix summarizes the properties of the distributions that we will need in this book.

1.5 *Summary of random variable theory

We can summarize the above informal concepts very compactly if we re-state them in mathematical form. A mathematical statement has the advantage not only of brevity but also of reducing ambiguity.

Formally, a random variable Y is defined as a function from a sample space of possible outcomes S to the real number system:

$$Y : S \rightarrow \mathbb{R} \quad (1.9)$$

The random variable associates to each outcome $\omega \in S$ exactly one number $Y(\omega) = y$. S_Y is all the y 's (all the possible values of Y , the support of Y). I.e., $y \in S_Y$.

Every random variable Y has associated with it a probability mass (distribution) function (PMF, PDF). I.e., PMF is used for discrete distributions and PDF for continuous distributions. The PMF/PDF maps every element of S_Y to a value between 0 and 1.

$$p_Y : S_Y \rightarrow [0, 1] \quad (1.10)$$

Probability mass functions (discrete case) and probability density functions (continuous case) are functions that assign probabilities or relative frequencies to all events in a sample space.

The expression

$$Y \sim f(\cdot) \quad (1.11)$$

will be used to mean that the random variable Y has pdf/pmf $f(\cdot)$. For example, if we say that $Y \sim \text{Binomial}(n, \theta)$, then we are asserting that the PMF is:

$$\text{Binomial}(k|n, \theta) = \binom{n}{k} \theta^k (1 - \theta)^{n-k} \quad (1.12)$$

If we say that $Y \sim \text{Normal}(\mu, \sigma)$, we are asserting that the PDF is

$$\text{Normal}(y|\mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y - \mu)^2}{2\sigma^2}\right) \quad (1.13)$$

The **cumulative distribution function** or CDF is defined as follows:

For discrete distributions, the probability that Y is less than a is written:

$$P(Y < a) = F(Y < a) = \sum_{-\infty}^a f(y) \quad (1.14)$$

For continuous distributions, the summation symbol \sum above becomes the summation symbol for the continuous case, which is the integral \int . The upper and lower bounds are marked by adding a subscript and a superscript on the integral. For example, if we want the area under the curve between points a and b for some function $f(y)$, we write $\int_b^a f(y) dy$. So, if we want the probability that Y is less than a , we would write:

$$P(Y < a) = F(Y < a) = \int_{-\infty}^a f(y) dy \quad (1.15)$$

The above integral is simply summing up the area under the curve between the points $-\infty$ and a ; this gives us the probability of observing a or a value smaller than a .

A final point here is that we can go back and forth between the PDF and the CDF. If the PDF is $f(y)$, then the CDF that allows us to compute quantities like $P(Y < b)$ is just the integral $F(Y < b) = \int_{-\infty}^b f(y) dy$. If we differentiate the CDF, we get the PDF back: $d(F(y))/dy = f(y)$.

1.6 Further reading

A quick review of the mathematical foundations needed for statistics is available in the short book by [Fox \(2009\)](#). [Morin \(2016\)](#) and [Blitzstein and Hwang \(2014\)](#) are accessible introductions to probability theory. [Kerns \(2010\)](#) contains a very well-written general introduction to random variable theory, but assumes the reader knows the basics of calculus.

1.7 Exercises

1.7.1 Practice using the `pnorm` function

1.7.1.1 Part 1

Given a normal distribution with mean 73 and standard deviation 101, use the `pnorm` function to calculate the probability of obtaining values between 30 and -47 from this distribution.

1.7.1.2 Part 2

Calculate the following probabilities. Given a normal distribution with mean 50 and standard deviation 3, what is the probability of getting

- a score of 40 or less
- a score of 40 or more
- a score of 59 or more

1.7.1.3 Part 3

Given a normal distribution with mean 55 and standard deviation 6, what is the probability of getting

- a score of 50 or less.
- a score between 52 and 58.
- a score of $\mu+1$ or more.

1.7.2 Practice using the `qnorm` function

1.7.2.1 Part 1

Consider a normal distribution with mean 1 and standard deviation 1.

Compute the lower and upper boundaries such that:

- the area (the probability) to the left of the lower boundary is 0.46.
- the area (the probability) to the left of the upper boundary is 0.81.

1.7.2.2 Part 2

Given a normal distribution with mean 59.061 and standard deviation 0.673. There exist two quantiles, the lower quantile q_1 and the upper quantile q_2 , that are equidistant from the mean 59.061, such that the area under the curve of the Normal probability between q_1 and q_2 is 85%. Find q_1 and q_2 .

1.7.3 Practice using qt

Take an independent random sample of size 157 from a normal distribution with mean 170, and standard deviation 37. Next, we are going to pretend we don't know the population parameters (the mean and standard deviation). We compute the MLEs of the mean and standard deviation using the data and get the sample mean 171.138 and the sample standard deviation 35.223.

- Compute the estimated standard error using the sample standard deviation provided above.
- What are your degrees of freedom for the relevant t-distribution?
- Calculate the **absolute** critical t-value for a 95% confidence interval using the relevant degrees of freedom you just wrote above.
- Next, compute the lower bound of the 95% confidence interval using the estimated standard error and the critical t-value.
- Finally, compute the upper bound of the 95% confidence interval using the estimated standard error and the critical t-value.

1.7.4 Maximum likelihood estimation 1

Given the data point 16.655. The function `dnorm` gives the likelihood given a data point (or multiple data points) and a value for the mean and the standard deviation (`sd`). Using `dnorm`, compute

- the likelihood of the data point 16.655 assuming a mean of 12 and standard deviation 5.
- the likelihood of the data point 16.655 assuming a mean of 11 and standard deviation 5.

- the likelihood of the data point 16.655 assuming a mean of 10 and standard deviation 5.
- the likelihood of the data point 16.655 assuming a mean of 9 and standard deviation 5.

1.7.5 Maximum likelihood estimation 2

You are given 10 independent and identically distributed data points that are assumed to come from a Normal distribution with unknown mean and unknown standard deviation:

```
x
```

```
## [1] 513 491 501 500 494 490 500 507 515 481
```

The function `dnorm` gives the likelihood given multiple data points and a value for the mean and the standard deviation. The log-likelihood can be computed by typing `dnorm(..., log=TRUE)`.

The product of the likelihoods for two independent data points can be computed like this: Suppose we have two independent and identically distributed data points 5 and 10. Then, assuming that the Normal distribution they come from has mean 10 and standard deviation 2, the joint likelihood of these is:

```
dnorm(5, mean = 10, sd = 2) * dnorm(10, mean = 10,
sd = 2)
```

```
## [1] 0.001748
```

It is easier to do this on the log scale, because then one can add instead of multiplying. This is because $\log(x \times y) = \log(x) + \log(y)$. For example:

```
log(2 * 3)
```

```
## [1] 1.792
```

```
log(2) + log(3)
```

```
## [1] 1.792
```

So the joint log likelihood of the two data points is:

```
dnorm(5, mean = 10, sd = 2, log = TRUE) + dnorm(10,
  mean = 10, sd = 2, log = TRUE)
```

```
## [1] -6.349
```

Even more compactly:

```
sum(dnorm(c(5, 10), mean = 10, sd = 2, log = TRUE))
```

```
## [1] -6.349
```

- Given the 10 data points above, calculate the maximum likelihood estimate (MLE) of the expectation.
- The sum of the log-likelihoods of the data-points x , using as the mean the MLE from the sample, and standard deviation 5.
- What is the sum of the log-likelihood if the mean used to compute the log-likelihood is 497.2?
- Which value for the mean, the MLE or 497.2, gives the higher log-likelihood?

1.7.6 Generating bivariate data

Generate 50 data points from two random variables X and Y , where $X \sim \text{Normal}(50, 100)$ and $Y \sim \text{Normal}(100, 20)$. The correlation between the random variables is 0.7. Plot the simulated data points from Y against those from X .

1.7.7 Generating multivariate data

The bivariate case can be generalized to more than two dimensions. Generate 50 data points from three random variables X , Y , and

Z , where $X \sim \text{Normal}(50, 100)$, $Y \sim \text{Normal}(100, 20)$, and $Z \sim \text{Normal}(200, 50)$. The correlation between the random variables X and Y is 0.5, between X and Z is 0.2, and between Y and Z is 0.7. Here, you will have to define a 3×3 variance covariance matrix, with the pairwise covariances in the off-diagonals. Plot the simulated data points as two-dimensional figures: Y against X , Y against Z , and X against Z .

2

Hypothetical repeated sampling and the t-test

This chapter introduces some of the foundational ideas behind hypothesis testing in the frequentist framework. The key idea is that of hypothetical repeated sampling. When we fit a model to a given data-set, we are assuming that this is one of potentially infinite numbers of exact repetitions of an experiment. We leverage some amazing properties of these exact repetitions in order to draw inferences from our particular data. The key idea to understand here is the central limit theorem, to which we turn next.

2.1 The central limit theorem using simulation

Suppose we collect some data, which can be represented by a vector y ; this is a *single sample*. Given data y , and assuming for concreteness that the underlying likelihood is a $Normal(\mu = 500, \sigma = 100)$, the sample mean and standard deviation, \bar{y} and s give us an estimate of the unknown parameters mean μ and the standard deviation σ that are assumed to generate the data. Figure 2.1 shows the distribution of a particular sample, where the number of data points is $n = 1000$. Note that in this example the parameters are specified by us, so they are not unknown; in a real data-collection situation, the sample mean and standard deviation are all we have as estimates of the parameters.

Suppose now that you had not a single sample of size 1000 but many repeated samples. This isn't something one can normally do in real life; we often run a single experiment or, at most, repeat the same experiment once. However, one can simulate repeated sampling easily within R. Let us take 100 repeated samples like

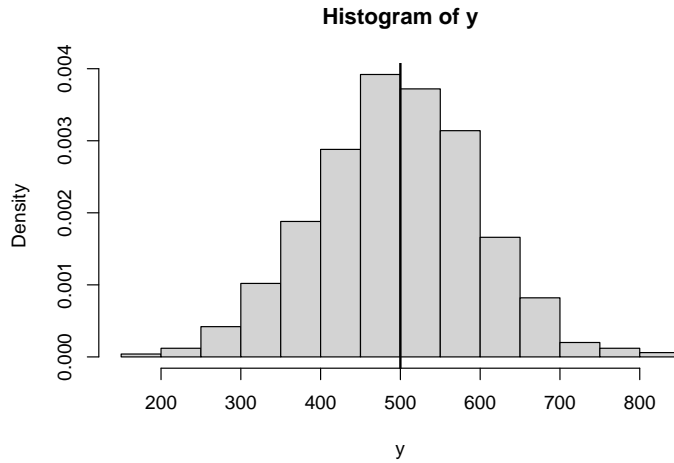


FIGURE 2.1: A sample of data y , from the distribution $\text{Normal}(500,100)$.

the one above, and save the samples in a matrix containing $n=1000$ rows and 100 columns, each column representing an experiment:

```
mu <- 500
sigma <- 100
## number of experiments:
k <- 100
## store for data:
y_matrix <- matrix(rep(NA, n * k), ncol = k)
for (i in 1:k) {
  ## expt result with sample size n:
  y_matrix[, i] <- rnorm(n, mean = mu, sd = sigma)
}
```

Now, if we compute the means \bar{y}_k of *each* of the $k = 1, \dots, 100$ experiments we just carried out, if certain conditions are met, these means will be normally distributed, with mean μ and standard deviation σ/\sqrt{n} . To understand this point, it is useful to first visualize the distribution of means and graphically summarize this standard deviation, which confusingly is called *standard error*.


```
## compute means from each replication:
y_means <- colMeans(y_matrix)
## the mean and sd (=standard error) of the means
mean(y_means)
```

```
## [1] 499.6
```

```
sd(y_means)
```

```
## [1] 3.053
```

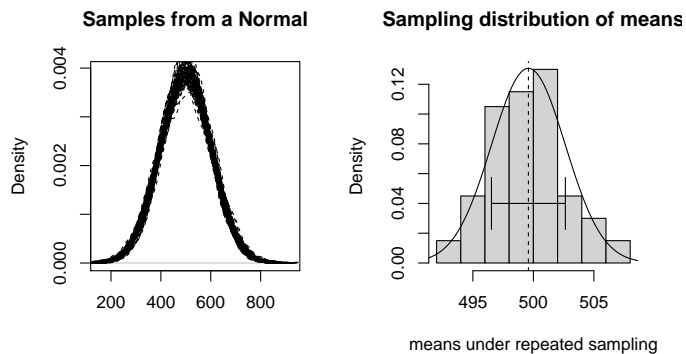


FIGURE 2.2: Sampling from a normal distribution (left); and the sampling distribution of the means under repeated sampling (right). The right-hand plot shows an overlaid normal distribution, and the standard deviation (standard error) as error bars.

The sampling distribution of means has a normal distribution provided two conditions are met: (a) the sample size should be large enough, and (b) μ and σ are defined for the probability density or mass function that generated the data. This fact is called the **central limit theorem** (CLT). The significance of the CLT for us as researchers is that from the summary statistics computed from a *single* sample, we can obtain an estimate of this distribution of means: $Normal(\bar{y}, s/\sqrt{n})$.

The statement that the sampling distribution of the means will be normal, with mean μ and standard deviation σ/\sqrt{n} , can be

derived formally through a surprisingly simple application of random variable theory. Suppose we gather independent and identically distributed data y_1, \dots, y_n , each of which is generated by a random variable $Y \sim \text{Normal}(\mu, \sigma)$.

When we compute the mean \bar{y} for each sample, we are assuming that each of the means is coming from a random variable \bar{Y} , which is just a linear combination of values generated by instances of the random variable Y , which itself has a pdf with mean (expectation) μ and variance σ^2 :

$$\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y = \frac{1}{n} Y_1 + \dots + \frac{1}{n} Y_n \quad (2.1)$$

So, the expectation of \bar{Y} is

$$\begin{aligned} E[\bar{Y}] &= E\left[\frac{1}{n} Y_1 + \dots + \frac{1}{n} Y_n\right] \\ &= \frac{1}{n} (E[Y] + \dots + E[Y]) \\ &= \frac{1}{n} (\mu + \dots + \mu) \\ &= \frac{1}{n} n\mu \\ &= \mu \end{aligned} \quad (2.2)$$

And the variance of \bar{Y} is

$$\begin{aligned} \text{Var}(\bar{Y}) &= \text{Var}\left(\frac{1}{n} Y_1 + \dots + \frac{1}{n} Y_n\right) \\ &= \frac{1}{n} [\text{Var}(Y_1 + \dots + Y_n)] \\ &= \frac{1}{n^2} \text{Var}(Y_1 + \dots + Y_n) \end{aligned} \quad (2.3)$$

The last line above arises because the variance of a random variable Z multiplied by a constant a , $\text{Var}(aZ)$ is $a^2 \text{Var}(Z)$. Here, $a = 1/n$, so $a^2 = 1/n^2$. Because Y_1, \dots, Y_n are independent, we can

compute the variance $Var(Y_1 + \dots + Y_n)$ by using the fact that the variance of the sum of independent random variables is the sum of their variances. This fact gives us:

$$\begin{aligned}
 \frac{1}{n^2} Var(Y_1 + \dots + Y_n) &= \frac{1}{n^2} (Var(Y) + \dots + Var(Y)) \\
 &= \frac{1}{n^2} n Var(Y) \\
 &= \frac{1}{n} Var(Y) \\
 &= \frac{\sigma^2}{n}
 \end{aligned} \tag{2.4}$$

This derives the above result that the expectation (i.e., the mean) and variance of the sampling distribution of the sample means are

$$E[\bar{Y}] = \mu \quad Var(\bar{Y}) = \frac{\sigma^2}{n} \tag{2.5}$$

The Central Limit Theorem, not proved here (for a proof, see p. 267 of [Miller and Miller \(2004\)](#)) can be summarized as follows.

Central Limit Theorem

Let $f(Y)$ be the pdf of a random variable Y , and assume that the pdf has mean μ and variance σ^2 . Then:

$$\bar{Y} \sim Normal(\mu, \sigma^2/n) \quad E[Y] = \mu, Var(Y) = \sigma^2 \quad \text{when } n \text{ is large} \tag{2.6}$$

For us, the practical implication of this result is huge. From a *single* sample y_1, \dots, y_n , we can derive the distribution of *hypothetical* sample means under repeated sampling. That is, it becomes possible to say something about what the plausible and implausible values of the sample mean are under repeated sampling. This is the basis for all hypothesis testing and statistical inference in the frequentist framework that we will look at in this book.

Sometimes the central limit theorem is misunderstood to imply

that the distribution that is assumed to generate the data is always going to be normal. It is important to understand that there are two distributions we are talking about here. First, there is the distribution that the data were generated from; this need not be normal. For example, you could get data from a Normal, Exponential, Gamma, or other distribution. Second, there is the sampling distribution of the *sample mean* under repeated sampling. It is the sampling distribution that the central limit theorem is about, not the distribution that generated the data.

2.2 Three examples of the sampling distribution

In the above discussion, the underlying pdf we sampled from above was a normal distribution. However, it need not be. Consider two examples: the underlying pdf is an Exponential or a Gamma distribution.

The Exponential distribution has a parameter λ (parameterized in R as a `rate`, $1/\lambda$); its mean is λ and its variance is $1/\lambda^2$. The sampling distribution is normal, even though the underlying distribution is an Exponential.

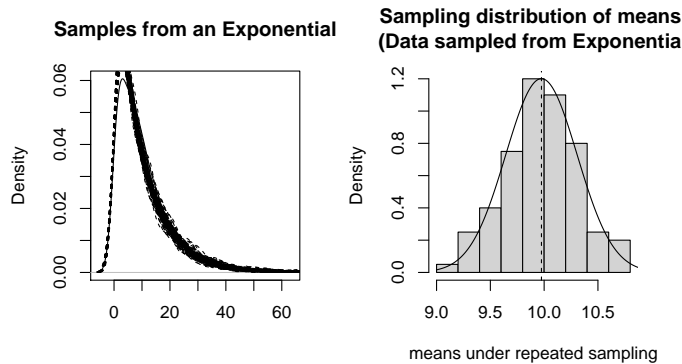


FIGURE 2.3: Sampling from an exponential.

A further example is samples from a Gamma distribution. Suppose we sample from a Gamma distribution with shape parameter cho-

sen arbitrarily to be 1. The distribution of means is again going to be approximately normal.

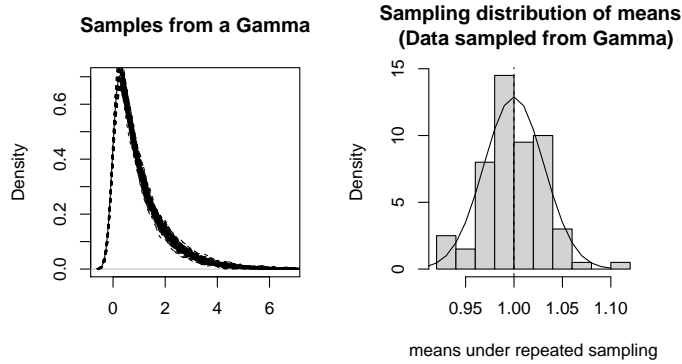


FIGURE 2.4: Sampling from a Gamma distribution.

As a final example, consider what happens if sample from a distribution, the Cauchy, that doesn't have any mean or variance defined for it.

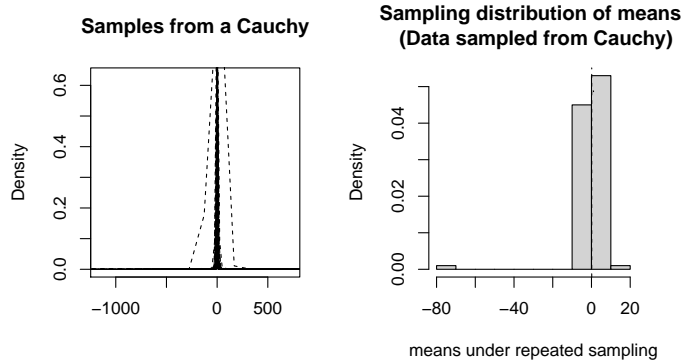


FIGURE 2.5: Sampling from a Cauchy distribution.

As the figure illustrates, when the mean and variance for the likelihood are undefined, the central limit theorem doesn't hold. In the rest of this book, we will always assume that the data are coming from a distribution that has a mean and variance defined for it.

2.3 The confidence interval, and what it's good for

Once we have sample of data y , and once the sample mean \bar{y} and the $SE = s/\sqrt{n}$ have been computed, it is common to define a so-called 95% confidence interval:

$$\bar{y} \pm 2SE \quad (2.7)$$

Because the sampling distribution of means is normally distributed, and because 95% of the area under the curve is covered by two times the standard deviation of the normal distribution, the upper and lower bounds of the interval defined by the interval $\bar{y} \pm 2SE$ covers approximately 95% of the area under the curve in the sampling distribution.

This interval is called the *confidence interval* (CI) and has the following meaning: If you take samples repeatedly and compute the CI each time, 95% of those CIs will contain the true population mean μ . One can simulate this situation. This time we will do 1000 repeated experiments instead of 100.

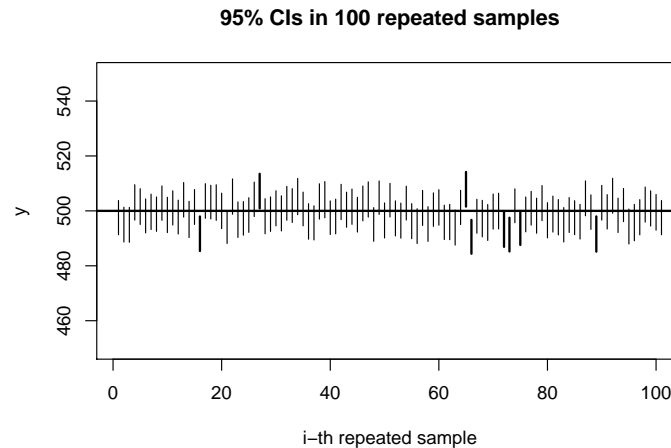
```
mu <- 500
sigma <- 100
n <- 1000
nsim <- 1000
lower <- rep(NA, nsim)
upper <- rep(NA, nsim)
for (i in 1:nsim) {
  y <- rnorm(n, mean = mu, sd = sigma)
  lower[i] <- mean(y) - 2 * sd(y)/sqrt(n)
  upper[i] <- mean(y) + 2 * sd(y)/sqrt(n)
}

## check how many CIs contain mu:
CIs <- ifelse(lower < mu & upper > mu, 1, 0)
```

```
## approx. 95% of the CIs contain true mean:
round(table(CIs)[2]/sum(table(CIs)), 2)
```

```
##      1
## 0.94
```

Figure 2.3 visualizes the coverage properties of the confidence interval in 100 simulations; by coverage we mean here the proportion of cases where the true μ is contained in the CI.



```
\begin{figure}
\caption{Illustration of the meaning of the 95% confidence inter-
val. The dark bars are CIs which do not contain the true mean.}
\end{figure}
```

The confidence interval is widely misinterpreted in a Bayesian way, i.e., as representing the range of plausible value of the μ parameter. This is the wrong interpretation because μ is a point value by assumption, it doesn't have a PDF associated with it. By contrast, the Bayesian credible interval does have this interpretation. In practice, we find that in most modeling settings we have encountered, the frequentist and Bayesian credible interval have very similar widths, with the Bayesian interval being slightly wider depending on the prior specifications.

Given the convoluted meaning of the CI, and the impossibility of interpreting a single CI, it is reasonable to ask: what good is a

CI? One can treat the CI as a graphical summary of width of the sampling distribution of the mean—the wider the sampling distribution, the more the implied variability under repeated sampling. The confidence interval can therefore be used as a visual summary of how uncertain we can be about the estimate of the sample mean under hypothetical repeated sampling. See [Cumming \(2014\)](#) for a useful perspective relating to using confidence intervals for inference.

We turn next to the central ideas behind the hypothesis test. We begin with the humble one-sample t-test, which contains many subtleties and is well worth close study before we move on to the main topic of this book: linear mixed models.

2.4 Hypothesis testing: The one sample t-test

With the central limit theorem and the idea of hypothetical repeated sampling behind us, we turn now to one of the simplest statistical tests that one can do with continuous data: the t-test.

Due to its simplicity, it is tempting to take only a cursory look at the t-test and move on immediately to the linear (mixed) model. This would be a mistake. The humble t-test is surprising in many ways, and holds several important lessons for us. There are subtleties in this test, and a close connection to the linear mixed model. For these reasons, it is worth slowing down and spending some time understanding this test. Once the t-test is clear, more complex statistical tests will be easier to follow, because the logic of these more complex tests will essentially be more of the same, or variations on this general theme. You will see later that t-test can be seen as an analysis of variance or ANOVA; and the paired t-test is exactly the linear mixed model with varying intercepts.

2.4.1 The one-sample t-test

As in our running example, suppose we have a random sample y of size n , and the data come from a $N(\mu, \sigma)$ distribution, with unknown parameters μ and σ . We can estimate μ from the sample mean \bar{y} , which we will sometimes also write as $\hat{\mu}$. We can also estimate σ from the sample standard deviation s , which we can also write as $\hat{\sigma}$. These estimates in turn allow us to estimate the sampling distribution of the mean under (hypothetical) repeated sampling:

$$N(\hat{\mu}, \frac{\hat{\sigma}}{\sqrt{n}}) \quad (2.8)$$

It is important to realize here that the above sampling distribution is only as realistic as the estimates of the mean and standard deviation parameters—if those are inaccurately estimated, then the sampling distribution is not realistic either.

Assume as before that we take an independent random sample of size 1000 from a random variable Y that is normally distributed, with mean 500 and standard deviation 100. As usual, begin by estimating the mean and SE:

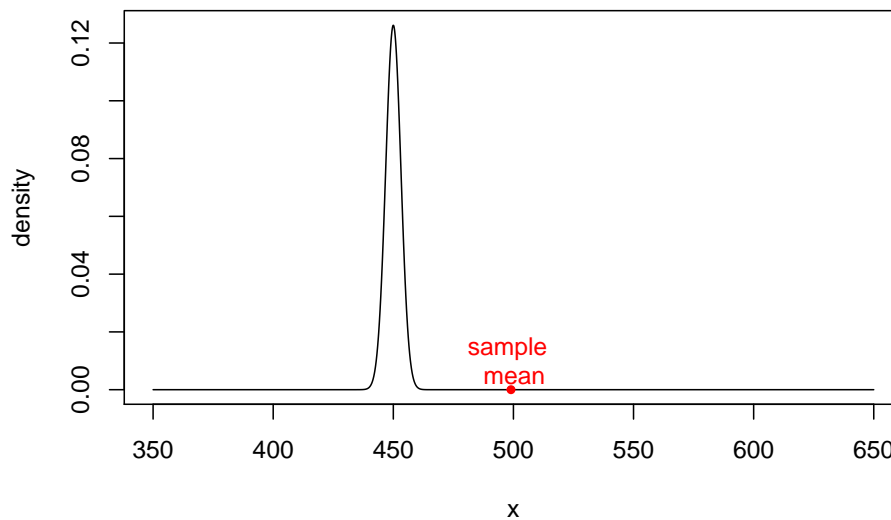
```
n <- 1000
mu <- 500
sigma <- 100
## generate simulated data:
y <- rnorm(1000, mean = 500, sd = 100)
## compute summary statistics:
y_bar <- mean(y)
SE <- sd(y)/sqrt(n)
```

The null hypothesis significance testing (NHST) approach as practised in psychology and other areas is to set up a null hypothesis that μ has some fixed value. Just as an example, assume that our null hypothesis is:

$$H_0 : \mu = 450 \quad (2.9)$$

This amounts to assuming that the true sampling distribution of sample means is (approximately) normally distributed and centered around 450, with the standard error estimated from the data.

The sampling distribution with $\mu=450$



The intuitive idea here is that - if the sample mean \bar{y} is “near” the hypothesized μ (here, 450), the data are possibly consistent with the null hypothesis distribution. - if the sample mean \bar{y} is “far” from the hypothesized μ , the data are inconsistent with the null hypothesis distribution.

The terms “near” and “far” will be quantified by determining how many standard error units the sample mean is from the hypothesized mean. This way of thinking shifts the focus away from the sampling distribution above, towards the distance measured in standard error units.

The distance between the sample mean and the hypothesized mean can be written in SE units. We will say that the sample mean is t standard errors away from the hypothesized mean:

$$t \times SE = \bar{x} - \mu \quad (2.10)$$

If we divide both sides with the standard error, we obtain the so-called observed t-statistic:

$$t = \frac{\bar{x} - \mu}{SE} \quad (2.11)$$

This observed t-value, an expression of the distance between the sample mean and the hypothesized mean, becomes the basis for the statistical test.

Notice that the t-value is a random variable: it is a transformation of \bar{X} , the random variable generating the sample means. The t-value can therefore be seen as an instance of the following transformed random variable T :

$$T = \frac{\bar{X} - \mu}{SE} \quad (2.12)$$

This random variable has a PDF associated with it, the t-distribution, which is defined in terms of the sample size n ; the pdf is written $t(n-1)$. Under repeated sampling, the t-distribution is generated from this random variable T .

We will compactly express the statement that “the observed t-value is assumed to be generated under repeated sampling from a t-distribution with $n-1$ degrees of freedom” as:

$$T \sim t(n-1) \quad (2.13)$$

For large n , the PDF of the random variable T approaches $N(0, 1)$. This is illustrated in Figure 2.6; notice that the t-distribution has fatter tails than the normal for small n , say $n < 20$, but for larger n , the t-distribution and the normal are essentially identical. Incidentally, when $n=2$, the t-distribution $t(1)$ is the Cauchy distribution we saw earlier; this distribution is characterized by fat tails, and has no mean or variance defined for it.

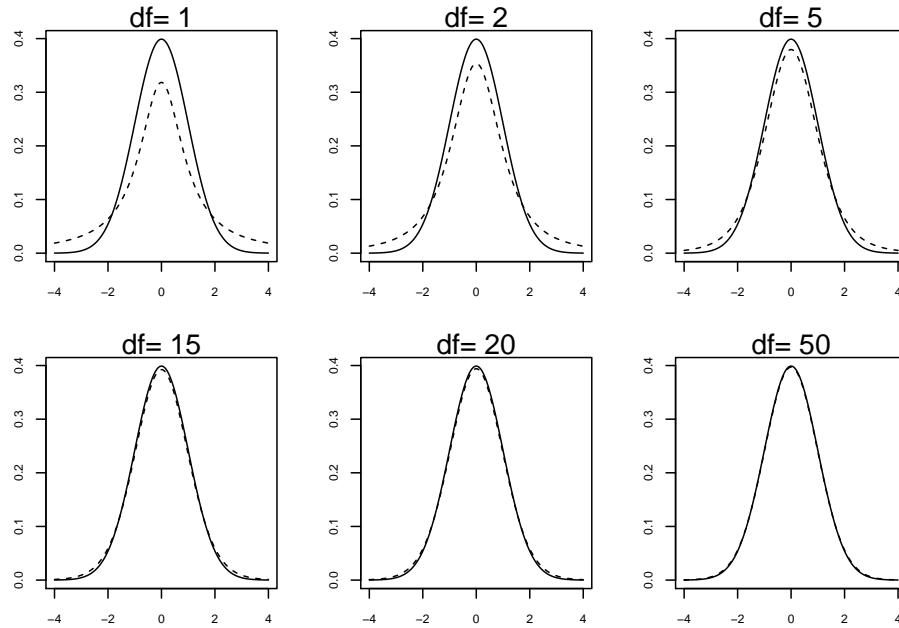


FIGURE 2.6: A visual comparison of the *t*-distribution (with degrees of freedom ranging from 1 to 50) with the standard normal distribution ($N(0,1)$).

Thus, given a sample size n , we can define a *t*-distribution corresponding to the null hypothesis distribution. For large values of n , we could even use $N(0,1)$, although it is traditional in psychology and linguistics to always use the *t*-distribution no matter how large n is.

The null hypothesis testing procedure proceeds as follows:

- Define the null hypothesis: in our example, the null hypothesis was that $\mu = 450$. This amounts to making a commitment about what fixed value we think the true underlying distribution of sample means is centered at.
- Given data of size n , estimate \bar{y} , standard deviation s , and from that, estimate the standard error s/\sqrt{n} . The standard error will be used to describe the sampling distribution's standard deviation.
- Compute the observed *t*-value:

$$t = \frac{\bar{y} - \mu}{s/\sqrt{n}} \quad (2.14)$$

- Reject null hypothesis if the observed t-value is “large” (to be made more precise next).
- Fail to reject the null hypothesis, or (under some conditions) even go so far as to accept the null hypothesis, if the observed t-value is “small”.

What constitutes a large or small observed t-value? Intuitively, the t-value from the sample is large when we end up far in *either* tail of the distribution. The two tails of the t-distribution will be referred to as the *rejection region*. The word *region* here refers to the real number line along the x-axis, under the tails of the distribution. The rejection region will go off to infinity on the outer sides, and is demarcated by a vertical line on the inner side of each tail. This is shown in Figure 2.7. It goes off to infinity because the support—the range of possible values—of the random variable that the t-distribution belongs to stretches from minus infinity to plus infinity.

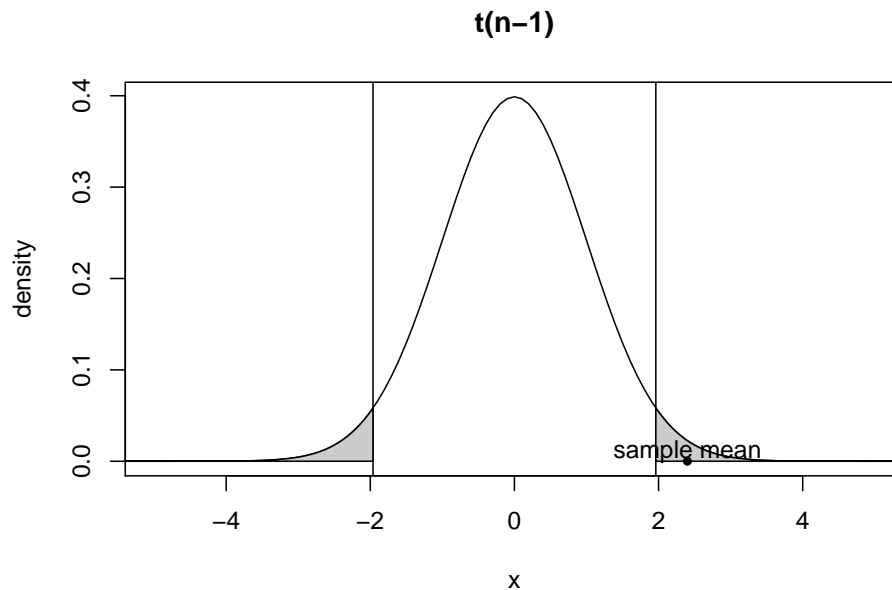


FIGURE 2.7: The rejection region in a t-distribution.

The location of the vertical lines is determined by the so-called *critical t-value* along the x-axis of the t-distribution. This is the value such that the area under the curve in the tails to the left or right of the tails is 0.025. As discussed in chapter 1, this area under the curve represents the probability of observing a value as extreme as the critical t-value, or some value that is more extreme. Notice that if we ask ourselves what the probability is of observing some particular t-value (a point value), the answer must necessarily be 0 (if you are unclear about why, re-read chapter 1). But we can ask the question: what is the absolute t-value, written $|t|$, such that $P(T > |t|) = 0.025$? That's the critical t-value.

For a given sample size n , we can identify the rejection region by using the `qt` function, whose usage is analogous to the `qnorm` function, discussed in chapter 1.

Because the shape of the t-distribution depends on the degrees of freedom ($n-1$), the critical t-value beyond which we reject the null hypothesis will change depending on sample size. For large sample sizes, say $n > 50$, the rejection point is about 2.

```
abs(qt(0.025, df = 15))
```

```
## [1] 2.131
```

```
abs(qt(0.025, df = 50))
```

```
## [1] 2.009
```

Consider the observed t-value from our sample in our running example:

```
## null hypothesis mean:
mu <- 450
(t_value <- (y_bar - mu)/SE)
```

```
## [1] 15.49
```

This observed t-value is huge and is telling you the distance of the sample mean from the null hypothesis mean μ in standard error units.

$$t = \frac{\bar{y} - \mu_0}{s/\sqrt{n}} \text{ or } t \frac{s}{\sqrt{n}} = \bar{y} - \mu_0 \quad (2.15)$$

For large sample sizes, if the absolute t-value $|t|$ is greater than 2 (approximately), we will reject the null hypothesis. The choice of 2 is purely conventional and comes from standard practice in psychology and related disciplines (as we will see in this book, standard practice is sometimes not a good-enough reason to decide on such things!).

For a smaller sample size n , you can compute the exact critical t-value:

```
qt(0.025, df = n - 1)
```

```
## [1] -1.962
```

Why is this critical t-value negative in sign? That is because it is on the left-hand side of the t-distribution, which is symmetric. The corresponding value on the right-hand side is:

```
qt(0.975, df = n - 1)
```

```
## [1] 1.962
```

These values are of course identical if we ignore the sign. This is why we always frame our discussion around the absolute t-value.

In R, the built-in function `t.test` delivers the observed t-value. Given our running example, with the null hypothesis $\mu = 450$, R returns the following:

```
## observed t-value from t-test function:
t.test(y, mu = 450)$statistic
```

```
##      t
## 15.49
```

The default value for the null hypothesis mean μ in this function is 0; so if one doesn't define a null hypothesis mean, the statistical test is done with reference to a null hypothesis that $\mu = 0$. That is why this t-value does not match our calculation above:

```
t.test(y)$statistic
```

```
##      t
## 157.8
```

In the most common usage of the t-test, the null hypothesis mean will be 0, because usually one is comparing a difference in means between two conditions or two sets of conditions. So the above line of code will work out correctly in those cases; but if you ever have a different null hypothesis mean than 0, then you have to specify it in the `t.test` function.

So, the t-test is used as if it furnishes a *decision rule*: either reject the null hypothesis or fail to reject it. Whenever we do an experiment and carry out a t-test, we use the t-test to make a decision: reject or fail to reject the null hypothesis.

Recall that behind the t-test lies the assumption that the observed t-value is coming from a random variable, $T \sim t(n - 1)$. The particular t-value we observe from a particular data-set belongs to a distribution of t-values under hypothetical repeated sampling. Thus, implicit in the logic of the t-test—and indeed every frequentist statistical test—is the assumption that the experiment is in principle repeatable: the experiment can in principle be re-run as many times as we want, assuming we have the necessary resources and time.

This implicit idea of the experiment's repeatability leads to an important aspect of the t-test: its long-run properties. In other words, its ability (at least in theory) to lead the researchers to the

correct decision in the long run, i.e., under (hypothetical) repeated sampling. We turn to this issue next.

2.4.2 Type I, II error, and power

When we do a hypothesis test using the t-test, the observed t-value will either fall in the rejection region, leading us to reject the null hypothesis, or it will land in the non-rejection region, leading us to fail to reject the null. That is a single, one-time event.

However, the null hypothesis can be either true or not true; we don't know which of those two possibilities is the reality. When we say that the null could be true, we mean that the parameter μ actually does have the hypothesized value μ_0 ; when we say that the null could be false, we mean that the parameter μ has some *specific* value μ_{alt} other than μ_0 . We can represent these two alternative possible realities in a tabular form, as shown below. The two columns show the two possible worlds, one in which the null is true, and the other in which it is false. The two rows show the two possible decisions we can take based on the observed t-value: reject the null or fail to reject it.

	Possible world 1	Possible world 2
	H_0 TRUE: $\mu = \mu_0$	H_0 FALSE $\mu = \mu_{alt}$
Decision: 'reject':	α	$1 - \beta$
	Type I error	Power
Decision: 'fail to reject':	$1 - \alpha$	β
		Type II error

As the table shows, we can make two kinds of mistakes:

- Type I error or α : Reject the null when it's true.
- Type II error or β : Accept the null when it's false.

In psychology and related areas, Type I error is fixed a priori at 0.05. This stipulated Type I error value is why the absolute critical t-value is kept at approximately 2; if, following recommendations from Benjamin et al. (2018), we were to stipulate that the Type I

error be 0.005, then the critical t-value would have had to be set at:

```
abs(qt(0.0025, df = n - 1))
```

```
## [1] 2.813
```

This suggested change in convention hasn't been taken up yet in cognitive science, but this could well change one day.

Type II error, the probability of incorrectly accepting the null hypothesis when it is false with some particular value for the parameter μ , is conventionally recommended (e.g., by the American Psychological Association) to be kept at 0.20 or lower. This implies that the probability of correctly rejecting a null hypothesis for some particular true value of μ is 1-Type II error. This probability, called statistical power, or just power, should then obviously be larger than 0.80. Again, there is nothing special about these stipulations; they are conventions that became the norm over time.

Next, we will consider the trade-off between Type I and II error. For simplicity, assume that the standard error is 1, and the null hypothesis is that $\mu = 0$. This means that the t-value is really the sample mean.

Consider the concrete situation where, in reality, the true value of μ is 2. As mentioned above, the null hypothesis H_0 is that $\mu = 0$. Now the H_0 is false because $\mu = 2$ and not 0. Type I and II error can be visualized graphically as shown in Figure 2.8.

To understand Figure 2.8, one has to consider two distributions side by side. First, consider the null hypothesis distribution, centered at 0. Under the null hypothesis distribution, the rejection region lies below the dark colored tails of the distributions. The area under the curve in these dark-colored tails is the Type I error (conventionally set at 0.05) that we decide on before we even conduct the experiment and collect the data. Because the Type I error is set at 0.05, and because the t-distribution is symmetric, the area under the curve in each tail is 0.025. The absolute critical

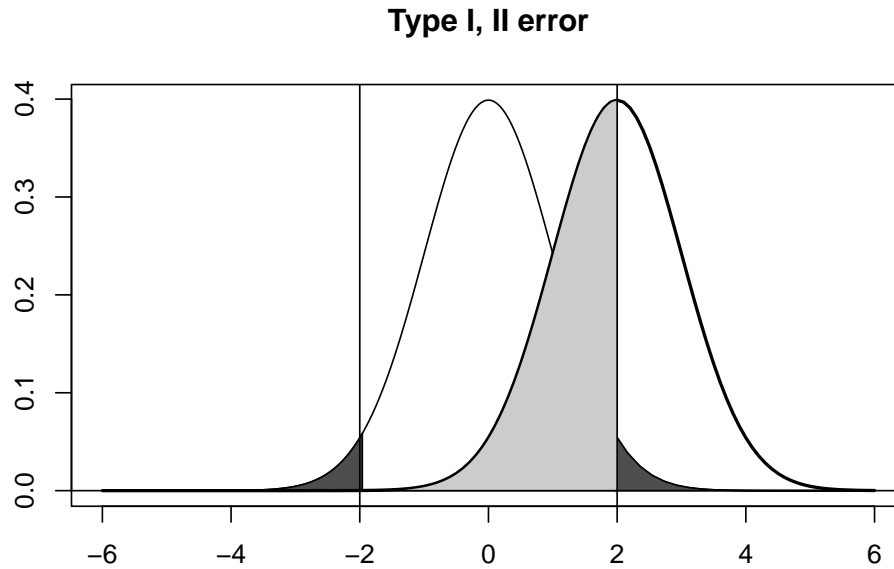


FIGURE 2.8: A visualization of Type I and II error. The dark-shaded tails of the left-hand side distribution represent Type I error; and the light-colored shaded region of the right-hand side distribution represents Type II error. Power is the unshaded area under the curve in the right-hand side distribution.

t-value helps us demarcates the inner boundaries of the rejection regions through the vertical lines shown in the figure. These vertical lines play a crucial role in helping us understand Type II error and power. This becomes clear when we consider the distribution representing the alternative possible value of μ , the distribution centered around 2. In this second distribution, consider now the area under the curve between the vertical lines demarcating the rejection region under the null. This area under the curve is the probability of accepting the null hypothesis when the null hypothesis is false with some specific value (here, when μ has value 2).

Some interesting observations follow. Suppose that the true effect is in fact $\mu = 2$, as in the above illustration. Then,

- Simply decreasing Type I error to a smaller value like 0.005 will also increase Type II error, which means that power (1-Type II error) will fall.

- Increasing sample size will squeeze the vertical lines closer to each other, reducing Type II error, and therefore increasing power. Decreasing sample size will have the opposite effect.
- If we design an experiment with a larger effect size, e.g., by setting up a stronger manipulation (concrete examples will be discussed in this book later on), our Type II error will go down, and therefore power will go up. Figure 2.9 shows a graphical visualization of a situation where the true effect size is $\mu = 4$. Here, Type II error is much smaller compared to Figure 2.8, where $\mu = 2$.

Type I, II error

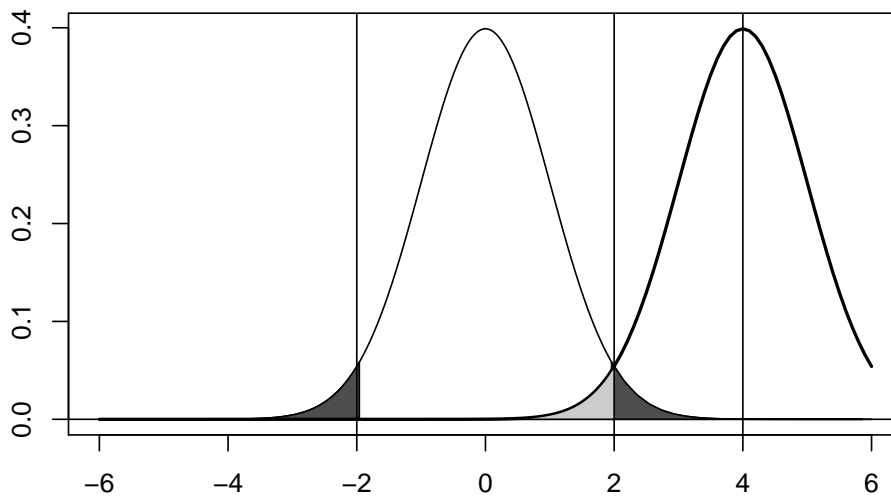


FIGURE 2.9: The change in Type II error if the true effect has mean 4.

In summary, when we plan out an experiment, we are also required to specify the Type I and II error associated with the design. Both sources of error are within our control, at least to some extent. The Type I error we decide to use will determine our critical *t*-value and therefore our decision criterion for rejecting, failing to reject, or even (under certain conditions, to be discussed below) accepting the null hypothesis.

The Type II error we decide on will determine the long-run prob-

ability of incorrectly accepting the null hypothesis; its inverse (1-Type II error), statistical power, will determine the long-run probability of correctly rejecting the null hypothesis under the assumption that the μ has some particular assumed value.

That's the theory anyway. In practice, researchers only rarely consider the power properties of their experiment design; the focus is almost exclusively on Type I error. The neglect of power in experiment design has had interesting consequences for theory development, as we will see later in this book. For a case study in psycholinguistics, see [Vasishth et al. \(2018\)](#).

2.4.3 How to compute power for the one-sample t-test

Power is a function of three variables:

- the effect size
- the standard deviation
- the sample size.

There are two ways that one can compute power in connection with the t-test: either one can use the built-in R function, `power.t.test`, or one can use simulation.

2.4.3.1 Power calculation using the `power.t.test`

Suppose that we have an expectation that an effect size is 15 ms \pm 5 ms (this could be based on the predictions of a theoretical model, or prior data); suppose also that prior experiments show standard deviations ranging from 100 to 300 ms. This is enough information to compute a power curve as a function of effect size and standard deviation. See Figure 2.10 and the associated code below.

```
sds <- seq(100, 300, by = 1)
lower <- power.t.test(delta = 15 - 5, sd = sds, n = 10,
  strict = TRUE)$power
upper <- power.t.test(delta = 15 + 5, sd = sds, n = 10,
  strict = TRUE)$power
```

```

meanval <- power.t.test(delta = 15, sd = sds, n = 10,
  strict = TRUE)$power

plot(sds, meanval, type = "l", main = "Power curve (n=10)\n using power.t.test",
  xlab = "standard deviation", ylab = "effect size")
lines(sds, lower, lty = 2)
lines(sds, upper, lty = 2)
text(200, 0.05, "10")
text(200, 0.054, "15")
text(200, 0.056, "20")

```

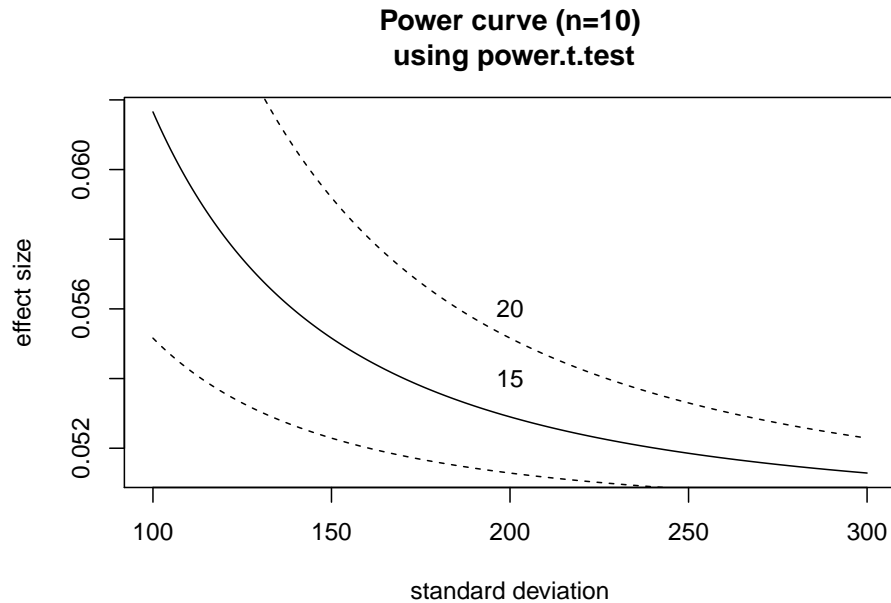


FIGURE 2.10: An illustration of a power curve for 10 participants, as a function of standard deviation, and three estimates of the effect: 15, 10, and 20.

2.4.3.2 Power calculations using simulation

An analogous calculation as the one shown above using the `power.t.test` function can also be done using simulated data. First, generate simulated data repeatedly for each possible com-

combination of parameter values (here, effect size and standard deviation), and compute the proportion of significant effects for each parameter combination. This can be done by defining a function that takes as input the number of simulations, sample size, effect size, and standard deviation:

```
compute_power <- function(nsim = 1000, n = 10, effect = NULL,
  stddev = NULL) {
  temp_power <- rep(NA, nsim)
  for (i in 1:nsim) {
    y <- rnorm(n, mean = effect, sd = stddev)
    temp_power[i] <- ifelse(abs(t.test(y)$statistic) >
      2, 1, 0)
  }
  ## return power calculation:
  mean(temp_power)
}
```

Then, plot the power curves as a function of effect size and standard deviation, exactly as in Figure 2.10. Power calculations using simulations are shown in Figure 2.11. It is clear that simulation-based power estimation is going to be noisy; this is because each time we are generating simulated data and then carrying out a statistical test on it. This is no longer a closed-form mathematical calculation as done in `power.t.test` (this function simply implements a formula for power calculation specified for this simple case). Because the power estimates will be noisy, we show a smoothed lowess line for each effect size estimate.

In the above example, simulation-based power calculation is overkill, and completely unnecessary because we have `power.t.test`. However, the technique shown above will be extended and will become our bread-and-butter method once we switch to power calculations for complicated linear mixed models. There, no closed form calculation can be done to compute power, at least not without oversimplifying the model; simulation will be the only practical way to calculate power.

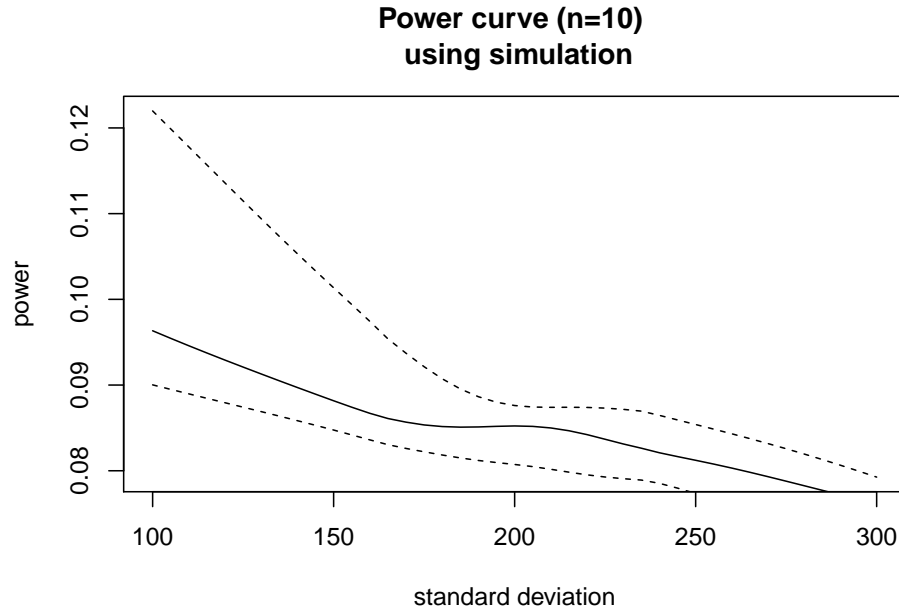


FIGURE 2.11: An illustration of a power curve using simulation, for 10 participants, as a function of standard deviation, and three estimates of the effect: 15, 10, and 20. The power curves are lowess-smoothed.

It is important to appreciate the fact that power is a *function*; it isn't a single number. Because we can never be sure what the true effect size is, or what the true standard deviation is, we should always aim to report a power function, as a function of plausible values of the relevant parameters.

2.4.4 The p-value

Continuing with our *t*-test example, the `t.test` function in R will not only print out a *t*-value as shown above, but also a probability known as a *p-value*. This is the probability of obtaining the observed *t*-value that we did, or some value more extreme than that, conditional on the assumption that the null hypothesis is true.

We can compute the *p*-value “by hand”. This can be computed, as done earlier, simply by calculating the area under the curve that

lies beyond the observed *t*-value. It is standard practice to take the tail probability on both sides of the *t*-distribution.

```
(t_value <- t.test(y, mu = 450)$statistic)
```

```
##      t  
## 15.49
```

```
2 * pt(abs(t_value), df = n - 1, lower.tail = FALSE)
```

```
##      t  
## 8.529e-08
```

The area from both sides of the tail is taken because it is conventional to do a so-called *two-sided t-test*: our null hypothesis is that $\mu = 450$, and our alternative hypothesis is two-sided: μ is either less than 450 or μ is larger than 450. When we reject the null hypothesis, we are accepting this alternative, that μ could be some value other than 450. Notice that this alternative hypothesis is remarkably vague; we would reject the null hypothesis regardless of whether the sample mean turns out to be 600 or -600, for example. The practical implication is that the *p*-value gives us the strength of the evidence against the null hypothesis; it doesn't give us evidence in favor of a specific alternative, such as saying that μ is positive or negative in sign. In psychology and allied disciplines, whenever the *p*-value falls below 0.05, it is common practice to write something along the lines that "there was reliable evidence for the predicted effect." This statement is technically incorrect; we only ever have evidence against the null. By looking at the sample mean and its sign, we are making a very big leap that we have evidence for the specific sample mean we happened to get. As we will see below, the sample mean can be wildly far from the true mean that produced the data.

One need not have done a two-sided alternative; one could have defined the alternative to be one-sided. In that case, one would compute only one side of the area under the curve. This kind of one-

sided test is not normally done, but one can imagine a situation where a one-sided test is justified (for example, when only one sign of the effect is possible, or if there is a strong theoretical reason to expect only one particular sign—positive or negative—on an effect). That said, in their scientific career, none of the authors of this book have ever had occasion to use a one-sided test.

The *p*-value is always interpreted with reference to the pre-defined Type I error. Conventionally, we reject the null if $p < 0.05$. This is because we set the Type I error at 0.05. Keep in mind that Type I error and the *p*-value are two distinct things. Type I error is the probability of your incorrectly rejecting the null under repeated sampling. This is not the same thing as your *p*-value. The latter will be obtained from a particular experiment, and will vary from experiment to experiment; it is a random variable. Type I error is a value we fix in advance.

2.4.4.1 *The distribution of the *p*-value under the null hypothesis

We have been talking about a continuous random variable as a dependent measure, and have learnt about the standard two-sided *t*-test, with a point null hypothesis. When we do such a test, we usually use the *p*-value to decide whether to reject the null hypothesis or not.

Sometimes, you will hear statisticians (e.g., Andrew Gelman on his blog) say that the null hypothesis significance test is a specific random number generator. What does that sentence mean? We explain this point here.

It's worth briefly reflecting on the fact that the *p*-value is a random variable; call it Z . The *p*-value is the cumulative distribution function (CDF) of the random variable T , which itself is a transformation of the random variable \bar{Y} :

$$T = (\bar{X} - \mu)/(\sigma/\sqrt{n})$$

This random variable T has some CDF $F(T)$. It is possible to show that if a random variable $Z = F(T)$, i.e., if Z is the CDF for the random variable T , then $Z \sim \text{Uniform}(0, 1)$.

This is an amazing fact. To get a grip on this, let's first think about the fact that when a random variable Z comes from a $Uniform(0, 1)$ distribution, then $P(Z < z) = z$. Consider some examples:

- when $z = 0$, then $P(Z < 0) = 0$;
- when $z = 0.25$, then $P(Z < 0.25) = 0.25$;
- when $z = 0.5$, then $P(Z < 0.5) = 0.5$;
- when $z = 0.75$, then $P(Z < 0.75) = 0.75$;
- when $z = 1$, then $P(Z < 1) = 1$.

Next, we will prove the above statement, that if a random variable $Z = F(T)$, i.e., if Z is the CDF for a random variable T , then $Z \sim Uniform(0, 1)$. The proof is actually quite astonishing and even has a name: it's called the *probability integral transform*.

Suppose that Z is the CDF of a random variable T : $Z = F(T)$. Then, it follows that $P(Z \leq z)$ can be rewritten in terms of the CDF of T : $P(F(T) \leq z)$. Now, if we apply the inverse of the CDF (F^{-1}) to both the left and right sides of the inequality, we get $P(F^{-1}F(T) \leq F^{-1}(z))$. But $F^{-1}F(T)$ gives us back T ; this holds because if we have a one-to-one onto function $f(x)$, then applying the inverse f^{-1} to this function gives us back x .

The fact that $F^{-1}F(T)$ gives us back T means that we can rewrite $P(F^{-1}F(T) \leq F^{-1}(z))$ as $P(T \leq F^{-1}(z))$. But this probability is simply the CDF $F(F^{-1}(z))$, which simplifies to z . This shows that $P(Z \leq z) = z$; i.e., that the p-value has a uniform distribution under the null hypothesis.

The above proof is restated below compactly:

$$\begin{aligned}
 P(Z \leq z) &= P(F(T) \leq z) \\
 &= P(F^{-1}F(T) \leq F^{-1}(z)) \\
 &= P(T \leq F^{-1}(z)) \\
 &= F(F^{-1}(z)) \\
 &= z
 \end{aligned} \tag{2.16}$$

It is for this reason that statisticians like Andrew Gelman peri-

odically point out that “the null hypothesis significance test is a specific random number generator”. The practical implication of this observation is that we should not place our theory development exclusively at the feet of the *p*-value. As we discuss in this book, other considerations (such as replicability, uncertainty of the estimates, and power) are as or even more important.

2.4.5 Type M and S error in the face of low power

Beyond Type I and II error, there are also two other kinds of error to be aware of. These are Type M and S error; both sources of error are closely related to statistical power.

The terms Type M and S error were introduced by [Gelman et al. \(2014\)](#), but the ideas has been in existence for some time ([Hedges, 1984](#), [Lane and Dunlap \(1978\)](#)). [Button et al. \(2013\)](#) refer to Type M and S error as the “winner’s curse” and “the vibration of effects.” In related work, [Ioannidis \(2008\)](#) discusses refers to the vibration ratio in the context of epidemiology.

Type S and M error can be illustrated with the following example. Suppose your true effect size is believed to be $D = 15$, then we can compute (apart from statistical power) the following error rates, which are defined as follows:

- **Type S error:** the probability that the sign of the effect is incorrect, given that the result is statistically significant.
- **Type M error:** the expectation of the ratio of the absolute magnitude of the effect to the hypothesized true effect size, given that result is significant. Gelman and Carlin also call this the exaggeration ratio, which is perhaps more descriptive than “Type M error”.

Suppose that our particular study has standard error 46, and sample size 37. And suppose that our true $\mu = 15$, as stated above. Then, we can compute statistical power, Type S and M error through simulation in the following manner:

```

## probable effect size, derived from past studies:
D <- 15
## SE from the study of interest:
se <- 46
stddev <- se * sqrt(37)
nsim <- 10000
drep <- rep(NA, nsim)
for (i in 1:nsim) {
  samp <- rnorm(37, mean = D, sd = stddev)
  drep[i] <- mean(samp)
}

```

Power can be computed by simply determining the proportion of times that the absolute observed t-value is larger than 2:

```

## power: the proportion of cases where we reject
## the null hypothesis correctly:
(pow <- mean(ifelse(abs(drep/se) > 2, 1, 0)))

```

```
## [1] 0.0608
```

Power is quite low here (we deliberately chose an example with low power to illustrate Type S and M error).

Next, we figure out which of the cases are statistically significant (which simulated values yield $p < 0.05$). As a criterion, we use a t-value of 2 to declare $p < 0.05$; we could have done this more precisely by working out an exact critical t-value.

```

## which results in drep are significant at
## alpha=0.05?
signif <- which(abs(drep/se) > 2)

```

Type S error is the proportion of significant cases with the wrong sign (sign error), and Type M error is the ratio by the true effect

(of $\mu = 15$) is exaggerated in those simulations that happened to come out significant.

```
## Type S error rate / signif:
(types_sig <- mean(drep[signif] < 0))
```

```
## [1] 0.1645
```

```
## Type M error rate / signif:
(typem_sig <- mean(abs(drep[signif])/D))
```

```
## [1] 7.414
```

In this scenario, when power is approximately 6%, whenever we get a significant effect, the probability of obtaining the wrong sign is a whopping 16% and the effect is likely to be 7.4144 times larger than its true magnitude. The practical implication is as follows.

When power is low, relying on the p-value (statistical significance) to declare an effect as being present will be misleading because the decision will be based on an overestimate of the effect (Type M error), and even the sign of the effect could be wrong. This isn't just a theoretical point; it has real-world consequences for theory development. For an example from psycholinguistics regarding this point, see [Vasishth et al. \(2018\)](#).

Another useful way to visualize Type M and S error is through the so-called funnel plot. As shown in Figure 2.12, estimates obtained from low-powered studies will tend to be exaggerated (the lower part of the funnel), and as power goes up, the effect estimates start to cluster tightly around the true value of the effect.

What is important to appreciate here is the fact that significant effects “point to the truth” just in case power is high; when power is low, either null results will frequently be found even if the null is false, and those results that turn out significant will be based on Type M error.

In many fields, it is practically impossible to conduct a high-

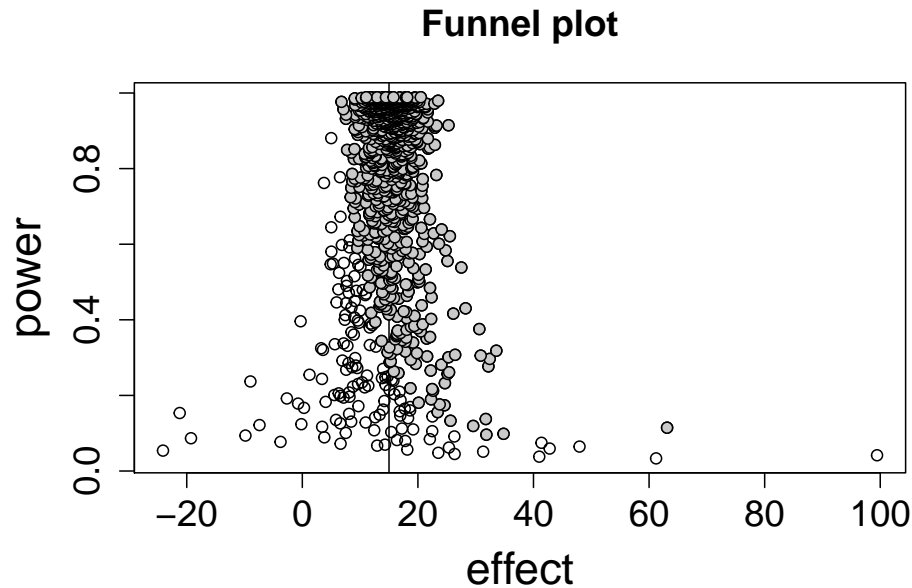


FIGURE 2.12: An illustration of a funnel plot. Shown are repeated samples of an effect estimate under different values of power, where the true value of the effect is 15 (marked by the vertical line). Significant effects are shaded gray. The lower the power, the wider the fluctuation of the effect; under low power, it is the exaggerated effects that end up statistically significant, even though they are very biased relative to the true value. As power goes up, the effect estimates start to cluster around the true value, and significant effects are also accurate estimates of the effect. Thus, low power leads to exaggerated estimates of the effect, especially if the data are filtered by statistical significance.

powered study. What should one do in this situation? When reporting results that are likely based on an underpowered study, the best approach is to openly acknowledge the power limitation, to attempt to conduct a direct replication of the effect to establish robustness, and to attempt to synthesize the evidence from existing knowledge ([Cumming, 2014](#)).

By direct replication, we mean that the study should be run multiple times with the same materials and design but new participants, to establish whether effect estimates in the original study and the

replication study are consistent with each other. Direct replications stand in contrast to so-called conceptual replications, which are not exact repetitions of the original design, but involve some further or slightly different but related experimental manipulation. Conceptual replications are also a very useful tool for cross-validating the existence of an effect.

Direct replications will always differ from the original study in some way or another—the lab may differ, the protocols might differ slightly, the experimenter is different, etc. Such between-study variability is obviously unavoidable in direct-replication attempts, but they are still worthwhile for establishing the existence of an effect. To make the idea of establishing robustness through replication attempts, detailed examples of different kinds of replication attempts will be presented in this book's example data-sets.

2.4.6 Searching for significance

The NHST procedure is essentially a decision procedure: if $p < 0.05$, we reject the null hypothesis; otherwise, we fail to reject the null. Because significant results are easier to publish than non-significant results, a common approach taken by researchers (including the first author of this book, when he was a graduate student) is to run the experiment and periodically check if statistical significance has been reached. The procedure can be described as follows:

- The experimenter gathers n data points, then checks for significance (is $p < 0.05$ or not?).
- If the result is not significant, he gets more data (say, n more data points). Then he checks for significance, and repeats.

Since time and money are limited, he might decide to stop collecting data after some multiple of n have been collected.

One can simulate different scenarios here. Suppose that n is initially 15.

Under the standard assumptions, we set Type I error to be 0.05. Let's suppose that the null hypothesis that $\mu = 0$ is in fact true, and that standard deviation is 250.


```

## Standard properties of the t-test:
pvals <- NULL
tstat_standard <- NULL
n <- 15
nsim <- 10000
## assume a standard dev of 1:
stddev <- 250
mn <- 0
for (i in 1:nsim) {
  samp <- rnorm(n, mean = mn, sd = stddev)
  pvals[i] <- t.test(samp)$p.value
  tstat_standard[i] <- t.test(samp)$statistic
}

```

Type I error rate is about 5%, consistent with our expectations:

```
round(table(pvals < 0.05)[2]/nsim, 2)
```

```
## TRUE
## 0.05
```

But the situation quickly deteriorates as soon as we adopt the strategy outlined above. Below, we will also track the distribution of the t-statistic.

```

pvals <- NULL
tstat <- NULL
## how many subjects can I run?
upper_bound <- n * 6

for (i in 1:nsim) {
  significant <- FALSE
  x <- rnorm(n, mean = mn, sd = stddev) ## take sample
  while (!significant & length(x) < upper_bound) {
    ## if not significant:

```

```

    if (t.test(x)$p.value > 0.05) {
      x <- append(x, rnorm(n, mean = mn, sd = stddev)) ## get more data
    } else {
      significant <- TRUE
    } ## otherwise stop:
  }
  pvals[i] <- t.test(x)$p.value
  tstat[i] <- t.test(x)$statistic
}

```

Now, Type I error rate is much higher than 5%:

```
round(table(pvals < 0.05)[2]/nsim, 2)
```

```
## TRUE
```

```
## 0.15
```

Figure 2.13 shows the distributions of the t-statistic in the standard case vs with the above stopping rule:

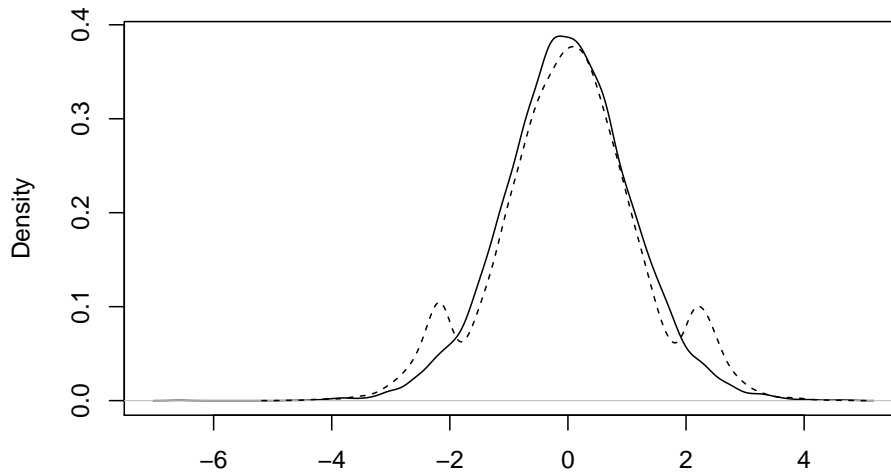


FIGURE 2.13: A comparison of the distribution of t-values with an a priori fixed stopping rule, versus a flexible stopping rule conditional on finding significance.

What is important to realize here is that the inflation in Type I error we observed above was due to the fact that the t-distribution is no longer a t-distribution: we have bumps in the tails when we use the flexible stopping rule, and these raise our Type I error. This demonstrates why one should fix one's sample size in advance, based on a power analysis. One should not deploy a stopping rule like the one above; if we used such a stopping rule, we are much more likely to incorrectly declare a result as statistically significant. There can be compelling reasons to adopt the peek-and-run strategy; e.g., if one wants to avoid exposing patients to a treatment that might turn out to be harmful. In such situations, one can run an adaptive experimental trial by correcting for Type I error inflation ([Pocock, 2013](#)). In this book, we will aim to develop a workflow whereby the sample size is fixed through power analysis, in advance of running an experiment.

2.5 The two-sample t-test vs. the paired t-test

In our running example above, we examined the case where we have a single vector of data y . This led to the one-sample t-test.

Next, we consider a case where we have two vectors of data. The data-set below is from [Johnson \(2011\)](#). Shown below are F1 formant data (in Herz), productions of different vowels by male and female speakers of different languages.

```
F1data <- read.table("data/F1_data.txt", header = TRUE)
F1data
```

```
##      female male vowel  language
## 1      391  339     i   W.Apache
## 2      561  512     e   W.Apache
## 3      826  670     a   W.Apache
## 4      453  427     o   W.Apache
## 5      358  291     i  CAEnglish
```

## 6	454	406	e	CAEnglish
## 7	991	706	a	CAEnglish
## 8	561	439	o	CAEnglish
## 9	398	324	u	CAEnglish
## 10	334	307	i	Ndumbea
## 11	444	361	e	Ndumbea
## 12	796	678	a	Ndumbea
## 13	542	474	o	Ndumbea
## 14	333	311	u	Ndumbea
## 15	343	293	i	Sele
## 16	520	363	e	Sele
## 17	989	809	a	Sele
## 18	507	367	o	Sele
## 19	357	300	u	Sele

Notice that the male and female values can be seen as *dependent* or *paired*: each row belongs to the same vowel and language. Nevertheless, we can compare males' and females' F1 frequencies, completely ignoring this paired nature of the data. The t-test does not “know” whether these data are paired or not—it is the researcher's job to make sure that model assumptions are met.

Let's ignore the paired nature of the data for now, and treat the two vectors as independent vectors. Suppose that our null hypothesis is that there is no difference between the mean F1's for males (μ_m) and females (μ_f). Now, our null hypothesis is $H_0 : \mu_m = \mu_f$ or $H_0 : \mu_m - \mu_f = \delta = 0$.

This kind of statistical test calls for a two-sample t-test. Note here that we are assuming that both the male and female F1 scores have equal variance.

The function call in R for a two-sample t-test is shown below:

```
t.test(F1data$female, F1data$male, paired = FALSE,
       var.equal = TRUE)
```

```
##
## Two Sample t-test
```

```
##
## data:  F1data$female and F1data$male
## t = 1.5, df = 36, p-value = 0.1
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -30.07 217.54
## sample estimates:
## mean of x mean of y
##      534.6      440.9
```

This t-test is computing the following t-statistic:

$$t = \frac{d - (\mu_m - \mu_f)}{SE} = \frac{d - 0}{SE} \quad (2.17)$$

where d is the difference between the two sample means; the rest of the terms we are familiar with. SE is the standard error of the sampling distribution of the difference between the means.

We will now do this calculation “by hand”. The only new things are the formula for the SE calculation, and the degrees of freedom for t-distribution $(2 \times n - 2) = 36$.

The standard error for the difference in the means in the two-sample t-test is computed using this formula:

$$SE_{\delta} = \sqrt{\frac{\hat{\sigma}_m^2}{n_m} + \frac{\hat{\sigma}_f^2}{n_f}} \quad (2.18)$$

Here, $\hat{\sigma}_m$ is estimate of the standard deviation for males, and $\hat{\sigma}_f$ for the females; the n are the respective sample sizes.

```
n_m <- n_f <- 19
## difference of sample means:
d <- mean(F1data$female) - mean(F1data$male)
(SE <- sqrt(var(F1data$male)/n_m + var(F1data$female)/n_f))

## [1] 61.04
```

```
(observed_t <- (d - 0)/SE)
```

```
## [1] 1.536
```

```
## p-value:
2 * (1 - pt(observed_t, df = 36))
```

```
## [1] 0.1334
```

The output of the two-sample t-test and the hand-calculation above match up.

Now consider what will change once we take into account the fact that the data are paired. The two-sample t-test now becomes a so-called paired t-test.

For such paired data, the null hypothesis is as before: $H_0 : \delta = 0$. But since each row in the data-frame is paired (from the same vowel+language), we subtract the vector row-wise, and get a new *vector* d (not a single number d as in the two-sample t-test) with the row-wise differences. Then, we just do the familiar one-sample test we saw earlier:

```
d <- F1data$female - F1data$male
t.test(d)
```

```
##
## One Sample t-test
##
## data: d
## t = 6.1, df = 18, p-value = 9e-06
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
## 61.48 125.99
## sample estimates:
## mean of x
## 93.74
```

An alternative syntax for the paired t-test explicitly feeds the two paired vectors into the function, but one must explicitly specify that they are paired, otherwise the test is a two-sample (i.e., unpaired) t-test:

```
t.test(F1data$female, F1data$male, paired = TRUE)

##
## Paired t-test
##
## data: F1data$female and F1data$male
## t = 6.1, df = 18, p-value = 9e-06
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  61.48 125.99
## sample estimates:
## mean of the differences
##                93.74
```

Incidentally, notice that the p-value in the paired t-test is statistically significant, unlike the two-sample t-test above. The null hypothesis is the same in both tests, but the significance level leads to different conclusions.

Which analysis is correct, the two-sample t-test or the paired t-test? It all depends on your assumptions about what the data represent. If you consider the data paired, for the reasons given above, then a paired test is called for. If the pairing (the same vowel and language in each row) is unlikely to create a dependency between the two data points in a row (here, domain knowledge is required), we can treat this as unpaired data.

Next, we look at some perhaps subtle points about the paired t-test.

2.5.1 Common mistakes involving the t-test

The paired t-test assumes that each row in the data-frame is independent of the other rows. This implies that the data-frame cannot

have more than one row for a particular pair. In other words, the data-frame cannot have repeated measurements spread out across rows.

For example, doing a paired *t*-test on this hypothetical data-frame would be incorrect:

female	male	vowel	language
391	339	i	W.Apache
400	320	i	W.Apache
⋮	⋮	⋮	⋮

Why? Because the assumption is that each row is independent of the others. This assumption is violated here (this is assuming that repeating the vowel from the same language will lead to some commonalities between the two repetitions).

Consider another hypothetical example. In the table below, from subject 1 we see two data points each for condition a and for condition b.

condition a	condition b	subject	item
391	339	1	1
400	320	1	2
⋮	⋮	⋮	⋮

Here, we again have repeated measurements from subject 1. The independence assumption is violated.

How to proceed when we have repeated measurements from each subject or each item? The solution is to aggregate the data so that each subject (or item) has only *one* value for each condition.

This aggregation allows us to meet the independence assumption of the *t*-test, but it has a potentially huge drawback: it pretends we have one measurement from each subject for each condition. Later on we will learn how to analyze unaggregated data, but if we want to do a paired *t*-test, we have no choice but to aggregate the data in this way.

A fully worked example will make this clear. We have repeated measures data on subject versus object relative clauses in English. The data are from a self-paced reading study reported in [Grodner and Gibson \(2005\)](#), their experiment 1. A theoretical prediction is that in English, object relatives are harder to read than subject relatives, in the relative clause verb region. We want to test this prediction.

First, load the data containing reading times from the region of interest (the relative clause verb):

```
gg05e1 <- read.table("data/grodnergibsonE1crit.txt",
                     header=TRUE)

head(gg05e1)
```

```
##      subject item condition rawRT
## 6          1    1   objgap   320
## 19         1    2   subjgap   424
## 34         1    3   objgap   309
## 49         1    4   subjgap   274
## 68         1    5   objgap   333
## 80         1    6   subjgap   266
```

We have repeated measurements for each condition from the subjects, and from items. You can establish this by using the `xtabs` command. Notice that there are no missing data points:

```
t(xtabs(~subject + condition, gg05e1))

##           subject
## condition 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
##   objgap  8 8 8 8 8 8 8 8 8  8  8  8  8  8  8  8  8  8
##   subjgap  8 8 8 8 8 8 8 8 8  8  8  8  8  8  8  8  8  8
##           subject
## condition 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33
```

```
##   objgap   8  8  8  8  8  8  8  8  8  8  8  8  8  8  8
##   subjgap  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8
##           subject
## condition 34 35 36 37 38 39 40 41 42
##   objgap   8  8  8  8  8  8  8  8  8
##   subjgap  8  8  8  8  8  8  8  8  8
```

```
t(xtabs(~item + condition, gg05e1))
```

```
##           item
## condition  1  2  3  4  5  6  7  8  9 10 11 12 13 14 15
##   objgap  21 21 21 21 21 21 21 21 21 21 21 21 21 21
##   subjgap 21 21 21 21 21 21 21 21 21 21 21 21 21 21
##           item
## condition 16
##   objgap  21
##   subjgap 21
```

It is important to stress once more that it is the researcher's responsibility to make sure that the t-test's assumptions are met. For example, one could fit a two-sample t-test to the data as provided. The two-sample t-test can be implemented using the syntax shown below:

```
t.test(rawRT ~ condition, gg05e1)
```

```
##
##  Welch Two Sample t-test
##
## data:  rawRT by condition
## t = 3.8, df = 431, p-value = 2e-04
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##   48.98 155.59
## sample estimates:
## mean in group objgap mean in group subjgap
##           471.4           369.1
```

This t-test is incorrect for several reasons, but the most egregious error here is that the data are paired (each subject delivers data for both conditions), and that property of the data is being ignored.

Another common mistake is to do a paired t-test on the data as provided. Again, the `t.test` function will happily return a meaningless result:

```
t.test(rawRT ~ condition, paired = TRUE, gg05e1)

##
## Paired t-test
##
## data: rawRT by condition
## t = 4, df = 335, p-value = 8e-05
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  51.98 152.59
## sample estimates:
## mean of the differences
##                102.3
```

Here, the degrees of freedom indicate that we have fit the incorrect model. There are 42 subjects and 16 items, and the presentation of items to subjects uses a Latin square design (each subject sees only one condition per item). The 335 degrees of freedom come from $42 \times 8 = 336$ data points, minus one. Why do we say 42×8 and not 42×16 ? That is because each subject will return eight differences in reading time for each condition: each subject gives us eight subject-relative data points and eight object-relative data points.

For each of the 42 subjects, the t-test function internally creates a vector of eight data points of subject relatives and subtracts the vector of eight data points of object relatives. That is how we end up with $42 \times 8 = 336$ data points.

These 336 data points are assumed by the t-test to be independent of each other; but this cannot be the case because each subject

delivers eight data points for each condition; these are obviously dependent (correlated) because they come from the same subject.

What is needed is a *single* data-point for each subject and condition, and for each item and condition. In order to conduct the *t*-test, aggregation of the data by subjects and by items is necessary.

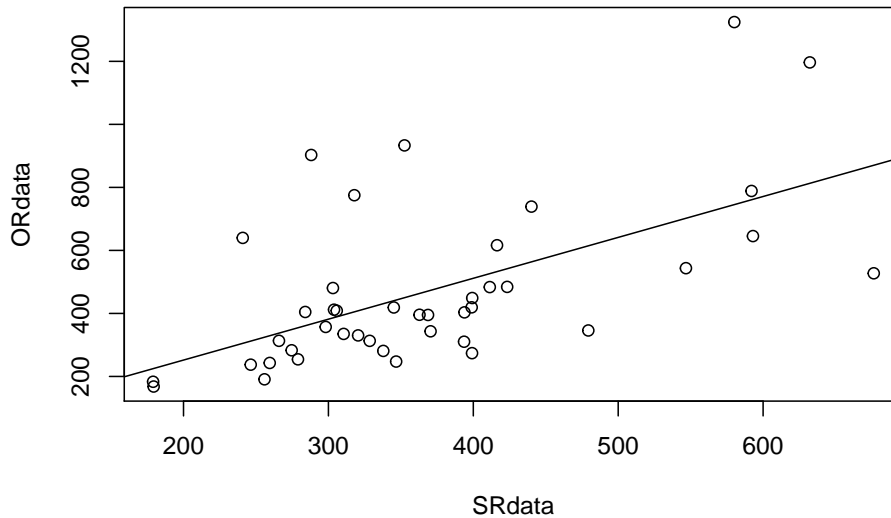
Consider the by-subjects aggregation procedure below. Now we have only one data-point for each condition and subject:

```
bysubj <- aggregate(rawRT ~ subject + condition, mean,
  data = gg05e1)
t(xtabs(~subject + condition, bysubj))
```

##		subject
##	condition	1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18
##	objgap	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##	subjgap	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##		subject
##	condition	19 20 21 22 23 24 25 26 27 28 29 30 31 32 33
##	objgap	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##	subjgap	1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##		subject
##	condition	34 35 36 37 38 39 40 41 42
##	objgap	1 1 1 1 1 1 1 1 1
##	subjgap	1 1 1 1 1 1 1 1 1

Notice that the data are correlated: the longer the subject relative clause data from a participant, the longer their object relative clause data:

```
SRdata <- subset(bysubj, condition == "subjgap")$rawRT
ORdata <- subset(bysubj, condition == "objgap")$rawRT
plot(SRdata, ORdata)
abline(lm(ORdata ~ SRdata))
```



```
cor(SRdata, ORdata)
```

```
## [1] 0.5876
```

Returning to the *t*-test, by aggregating the data the *t*-test's assumptions are met, and the degrees of freedom for this by-subjects analysis are now correct ($42 - 1 = 41$):

```
t.test(rawRT ~ condition, bysubj, paired = TRUE)
```

```
##
## Paired t-test
##
## data: rawRT by condition
## t = 3.1, df = 41, p-value = 0.003
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## 35.85 168.72
## sample estimates:
## mean of the differences
## 102.3
```

Similar to the by-subjects aggregation done above, one could do

a by-items aggregation and then a by-items t-test (What should be the degrees of freedom for the by-items analysis? There are 16 items in this data-set). This is left as an exercise for the reader.

The paired t-test illustrated above is actually not the best way to analyze this data-set, because it ignores the fact that each subject delivers not one but eight data points per condition. Each subject's repeated measurements will introduce a source of variance, but this source of variance is being suppressed in this t-test, leading to a possibly over-enthusiastic t-value. In order to take this variability into account, we must switch to the linear mixed model. But before we get to the linear mixed model, we have to consider the linear model. The next chapter turns to this topic.

2.6 Exercises

2.6.1 Computing the p-value

A paired t-test is done with data from 10 participants. The t-value from the test is 2.1. What is the p-value associated with a two-sided null hypothesis test?

2.6.2 Computing the t-value

If the p-value from a two-sided null hypothesis test had been 0.09, what would be the associated absolute t-value (i.e., ignoring the sign on the t-value)? The number of participants is 10, as above.

2.6.3 Type I and II error

Given that Type I error is 0.01; what is the highest value possible for Type II error?

3

Linear models and linear mixed models

3.1 From the t-test to the linear (mixed) model

We begin with the [Grodner and Gibson \(2005\)](#) self-paced reading data we saw in the previous chapter. Load the data and compute the means by condition at the relative clause verb:

```
gg05e1 <- read.table("data/grodnergibsonE1crit.txt",
  header = TRUE)
means <- round(with(gg05e1, tapply(rawRT, IND = condition,
  mean)))
means
```

```
## objgap subjgap
##      471      369
```

As predicted by theory, object relatives (labeled objgap here) are read slower than subject relatives (labeled subjgap).

As discussed in the previous chapter, a paired t-test can be done to evaluate whether we have evidence against the null hypothesis that object relatives and subject relatives have identical reading times. However, we have to aggregate the data by subjects and by items first.

```
bysubj <- aggregate(rawRT ~ subject + condition, mean,
  data = gg05e1)
byitem <- aggregate(rawRT ~ item + condition, mean,
```

```
data = gg05e1)
t.test(rawRT ~ condition, paired = TRUE, bysubj)$statistic

##      t
## 3.109
```

```
t.test(rawRT ~ condition, paired = TRUE, byitem)$statistic

##      t
## 3.754
```

What these two t-tests show is that both by subjects and by items, there is strong evidence against the null hypothesis that the object and relatives have identical reading times.

Interestingly, exactly the same t-values can be obtained by running the following commands, which implement a kind of linear model called the *linear mixed model*:

```
library(lme4)
```

```
## Loading required package: Matrix
```

```
m0lmersubj <- lmer(rawRT ~ condition + (1 | subject),
  bysubj)
summary(m0lmersubj)$coefficients
```

```
##              Estimate Std. Error t value
## (Intercept)      471.4       31.13  15.143
## conditionsubjgap -102.3       32.90  -3.109
```

```
m0lmeritem <- lmer(rawRT ~ condition + (1 | item),
  byitem)
summary(m0lmeritem)$coefficients
```


##	Estimate	Std. Error	t value
## (Intercept)	471.4	20.20	23.336
## conditionsubjgap	-102.3	27.25	-3.754

The signs of the t-values are the opposite to that of the paired t-tests above; the reason for that will presently become clear.

Our goal in this chapter is to understand the above model involving the `lmer` function, using the familiar paired t-test as a starting point.

For now, consider only the by-subject analysis. Given the sample means shown above for the two conditions, we can rewrite our best guess about how the object and subject relative clause reading time distributions were generated:

- Object relative: $Normal(471 - 102 \times 0, \hat{\sigma})$
- Subject relative: $Normal(471 - 102 \times 1, \hat{\sigma})$

Note that the two distributions for object and subject relative are assumed to be independent. This assumed independence is expressed by the fact that we define two separate Normal distributions, one for object relatives and the other for subject relatives. We saw earlier that this independence assumption of independence does not hold in our data, because we have one data point for each RC type from the same subject. However, for now we will ignore this detail; we will fix this shortcoming later.

The interesting point to notice here is that the mean for the object and subject relatives' distributions can be rewritten as a sum of two terms. A completely equivalent way to express the fact that object relatives are coming from a $Normal(471, \hat{\sigma})$ is to say that each object relative data-point can be described by the following equation:

$$y = 471 + -102 \times 0 + \varepsilon \text{ where } \varepsilon \sim Normal(0, \hat{\sigma}) \quad (3.1)$$

Similarly, the subject relative's distribution can be written as being generated from:

$$y = 471 - 102 \times 1 + \varepsilon \text{ where } \varepsilon \sim \text{Normal}(0, \hat{\sigma}) \quad (3.2)$$

In these data, the parameter $\hat{\sigma}$ is estimated to be 213. How do we know what this estimate is? This parameter's estimate can be derived from the by-subjects t-test output above: The observed t-value is

$$obs.t = \frac{\bar{x}}{s/\sqrt{n}} \quad (3.3)$$

Solving for s :

$$s = \bar{x} \times \sqrt{n}/obs.t = -103 \times \sqrt{42}/-3.109 = 213 \quad (3.4)$$

So, our model for the relative clause data consists of two equations:

Object relatives:

$$y = 471 - 102 \times 0 + \varepsilon \text{ where } \varepsilon \sim \text{Normal}(0, 213) \quad (3.5)$$

Subject relatives:

$$y = 471 - 102 \times 1 + \varepsilon \text{ where } \varepsilon \sim \text{Normal}(0, 213) \quad (3.6)$$

The above statements describe a *generative process* for the data.

Given such a statement about the generative process, we can express the estimated mean reading times for each RC type as follows. We can ignore the term ε because it has mean 0 (we stipulate this when we specify that $\varepsilon \sim \text{Normal}(0, \sigma)$).

Mean object relative reading times:

$$\text{Mean OR RT} = 471 - 102 \times 0 \quad (3.7)$$

Mean subject relative reading times:

$$\text{Mean SR RT} = 471 - 102 \times 1 \quad (3.8)$$

There is a function in R, the `lm()` function, which expresses the above statistical model, and prints out exactly the same numerical values that we used above:

```
summary(m0 <- lm(rawRT ~ condition, bysubj))$coefficients
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      471.4        31.13  15.143 1.795e-25
## conditionsubjgap -102.3        44.02  -2.324 2.263e-02
```

The linear model function `lm()` prints out two coefficients, 471 and -102 , that help express the mean reading times for object and subject relative data, using a simple coding scheme: object relatives are coded as 0, and subject relatives are coded as 1. This coding scheme is not visible to the user, but is represented internally in R. The user can see the coding for each condition level by typing:

```
## make sure that the condition column is of type
## factor:
bysubj$condition <- factor(bysubj$condition)
contrasts(bysubj$condition)
```

```
##      subjgap
## objgap      0
## subjgap      1
```

We will discuss coding in detail in a later chapter, but right now the simple 0,1 coding above—called treatment contrasts—is enough for our purposes.

Thus, what the linear model above gives us is two numbers: the mean object relative reading time (471), and the *difference* between object and subject relative (-102). We can extract the two coefficients by typing:

```
round(coef(m0))
```

```
##      (Intercept) conditionsubjgap
##                471                -102
```

In the vocabulary of linear modeling, the first number is called the *intercept*, and the second one is called the *slope*. Note that the meaning of the intercept and slope depends on the ordering of the factor levels. We can make the sample mean of the subject relative represent the intercept:

```
## reverse the factor level ordering:
bysubj$condition <- factor(bysubj$condition, levels = c("subjgap",
  "objgap"))
contrasts(bysubj$condition)
```

```
##      objgap
## subjgap    0
## objgap     1
```

Now, the intercept is the mean of the subject relatives, and the slope is the difference between object and subject relative reading time. Note that the sign of the t-value has changed—the sign depends on the contrast coding.

```
m1a <- lm(rawRT ~ condition, bysubj)
summary(m1a)$coefficients
```

```
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      369.1       31.13  11.857 1.819e-19
## conditionobjgap   102.3       44.02   2.324 2.263e-02
```

Let's switch back to the original factor level ordering:

```
bysubj$condition <- factor(bysubj$condition, levels = c("objgap",
  "subjgap"))
contrasts(bysubj$condition)
```

```
##      subjgap
```

```
## objgap      0
## subjgap     1
```

In mathematical form, the model can now be stated as a single equation:

$$rawRT = \beta_0 + \beta_1 condition + \varepsilon \quad (3.9)$$

where

- condition is a 0,1 coded vector, with object relatives coded as 0, and subject relatives coded as 1.
- β_0 is the mean for the object relative (which is coded as 0)
- β_1 is the amount by which the object relative mean must be changed to obtain the mean for the subject relative.
- ε is the noisy variation from trial to trial around the means for the two conditions, represented by $Normal(0, 213)$.

The null hypothesis of scientific interest here is always with reference to the slope, that the difference in means between the two relative clause types β_1 is:

$$H_0 : \beta_1 = 0$$

The t-test value printed out in the linear model is simply the familiar t-test formula in action:

$$obs.t = \frac{\beta_1 - 0}{SE} \quad (3.10)$$

The intercept also has a null hypothesis associated with it, namely that $H_0 : \beta_0 = 0$. However, this null hypothesis test is of absolutely no interest for us. This hypothesis test is reported by the `lm()` function only because the intercept is needed for technical reasons, to be discussed later.

The *contrast coding* mentioned above determines the meaning of the β parameters:

```
bysubj$condition <- factor(bysubj$condition, levels = c("objgap",
  "subjgap"))
contrasts(bysubj$condition)
```

```
##          subjgap
## objgap          0
## subjgap          1
```

When discussing linear models, we will make a distinction between the unknown true means β_0, β_1 and the estimated mean from the data $\hat{\beta}_0, \hat{\beta}_1$. The estimates that we have from the data are:

- Estimated mean object relative processing time: $\hat{\beta}_0 = 471$.
- Estimated mean subject relative processing time: $\hat{\beta}_0 + \hat{\beta}_1 = 471 + -102 = 369$.

3.2 Sum coding

We have established so far that the mathematical form of the model is:

$$rt = \beta_0 + \beta_1 condition + \varepsilon \quad (3.11)$$

We can change the contrast coding of the `condition` vector in the following way. First, recode the levels of the condition column as shown below.

```
## new contrast coding:
bysubj$cond <- ifelse(bysubj$condition == "objgap",
  1, -1)
```

Now, the two conditions are coded not as 0, 1 but as -1 and +1:

```
xtabs(~cond + condition, bysubj)
```

```
##      condition
## cond objgap subjgap
##   -1      0     42
##    1     42      0
```

With this coding, the model parameters have a different meaning:

```
m1 <- lm(rawRT ~ cond, bysubj)
round(coef(m1))
```

```
## (Intercept)      cond
##          420         51
```

- The intercept now represents the grand mean processing time: $\hat{\beta}_0 = 420$.
- The mean object relative processing time is now: $\hat{\beta}_0 + \hat{\beta}_1 \times 1 = 420 + 51 = 471$.
- The mean subject relative processing time is: $\hat{\beta}_0 + \hat{\beta}_1 \times (-1) = 420 - 51 = 369$.

This kind of parameterization is called *sum-to-zero contrast* or more simply *sum contrast* coding. This is the coding we will use most frequently in this book. We will elaborate on contrast coding in a later chapter; there, the advantages of sum coding over treatment coding will become clear. For now, it is sufficient to understand that one can *reparametrize* the model using different contrast codings, and that such a reparametrization impacts the interpretation of the parameters.

With sum coding, the null hypothesis for the slope is

$$H_0 : 1 \times \mu_{obj} + (-1) \times \mu_{subj} = 0 \quad (3.12)$$

The sum contrast coding of +1 standing for object relatives and -1 standing for subject relatives in the linear model directly refer to

the ± 1 coefficients in the null hypothesis above. Now the model is as follows.

Object relative reading times:

$$rt = 420 \times \mathbf{1} + 51 \times \mathbf{1} + \varepsilon \quad (3.13)$$

Subject relative reading times:

$$rt = 420 \times \mathbf{1} + 51 \times (-\mathbf{1}) + \varepsilon \quad (3.14)$$

One could write it in a single line as:

$$rt = 420 + 51 \times condition + \varepsilon \quad (3.15)$$

3.3 Checking model assumptions

It is an assumption of the linear model that the residuals are (approximately) normally distributed. That is what the statement $\varepsilon \sim Normal(0, \sigma)$ implies. It is important to check that model assumptions are approximately satisfied; this is because the null hypothesis significance testing procedure requires approximate normality of residuals.

Here is how we can check whether this normality assumption is met:

```
## extract residuals:
res.m1 <- residuals(m1)
```

Compare the residuals to the quantiles of the standard normal distribution ($Normal(0, 1)$):

```
library(car)
```



```
## Loading required package: carData

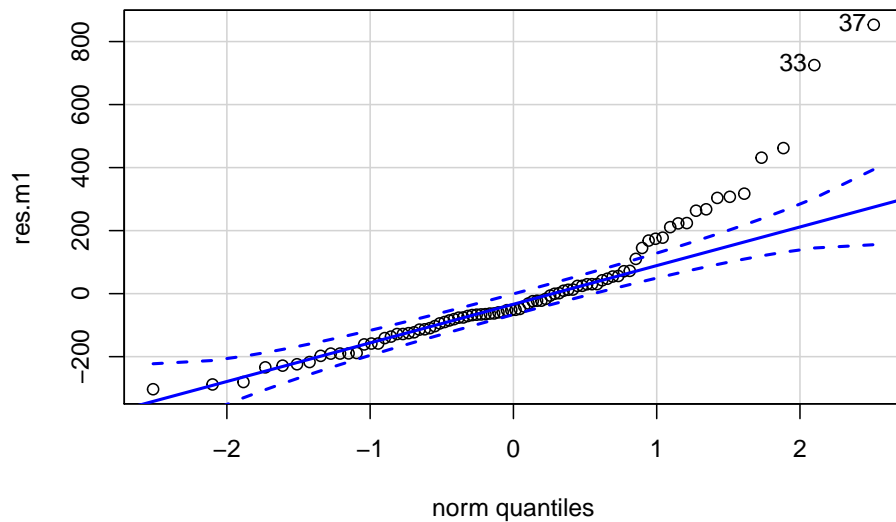
## Registered S3 methods overwritten by 'car':
##   method                      from
##   influence.merMod             lme4
##   cooks.distance.influence.merMod lme4
##   dfbeta.influence.merMod      lme4
##   dfbetas.influence.merMod     lme4

##
## Attaching package: 'car'

## The following object is masked from 'package:purrr':
##
##   some

## The following object is masked from 'package:dplyr':
##
##   recode
```

```
qqPlot(res.m1)
```



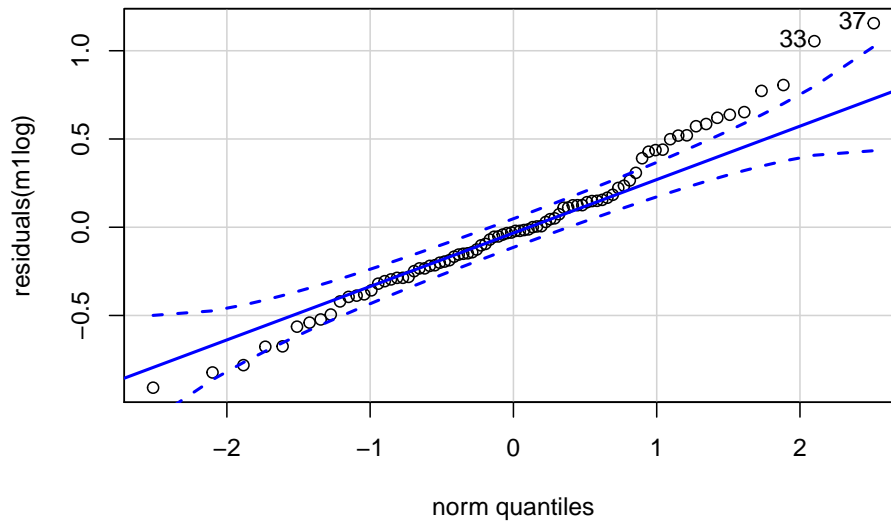
```
## [1] 37 33
```

When the normality assumption is met, the residuals will align

perfectly with the quantiles of the standard normal distribution, resulting in a straight diagonal line in the above plot. When the normality assumption is not met, the line will tend to curve away from the diagonal.

In the above case, a log transform of the data improves the normality of residuals:

```
m1log <- lm(log(rawRT) ~ cond, bysubj)
qqPlot(residuals(m1log))
```



```
## [1] 37 33
```

The estimates of the parameters are now in the log scale:

- The estimated grand mean processing time: $\hat{\beta}_0 = 5.9488$.
- The estimated mean object relative processing time: $\hat{\beta}_0 + \hat{\beta}_1 = 5.9488 + 0.0843 = 6.0331$.
- The estimated mean subject relative processing time: $\hat{\beta}_0 - \hat{\beta}_1 = 5.9488 - 0.0843 = 5.8645$.

The model does not change, only the scale does:

$$\log rt = \beta_0 + \beta_1 \text{condition} + \varepsilon \quad (3.16)$$

Now, the intercept and slope can be used to compute the reading time in the two conditions:

- Estimated mean object relative reading time: $\hat{\beta}_0 + \hat{\beta}_1 = 5.9488 + 0.0843 = 6.0331$.
- Estimated mean subject relative reading time: $\hat{\beta}_0 - \hat{\beta}_1 = 5.9488 - 0.0843 = 5.8645$.

Note that because $\exp(\log(rt)) = rt$, to get the mean estimates on the raw ms scale, we just need to exponentiate both sides of the equation:

$$\exp(\log rt) = \exp(\beta_0 + \beta_1 \text{condition}) \quad (3.17)$$

This approach gives us the following estimates on the ms scale:

- Estimated mean object relative reading time: $\exp(\hat{\beta}_0 + \hat{\beta}_1) = \exp(5.9488 + 0.0843) = 417$.
- Estimated mean subject relative reading time: $\exp(\hat{\beta}_0 - \hat{\beta}_1) = \exp(5.9488 - 0.0843) = 352$.

The difference in reading time is $417 - 352 = 65$ ms (cf. 102 ms from the model using the raw scale). The larger estimate based on the raw scale is less realistic, and we will see later that the large difference between the two conditions is driven by a few extreme, influential values.

In a later chapter, we will return to the issue of which transform to choose, and what the options are that are available for different situations.

3.4 From the paired t-test to the linear mixed model

One interesting point to notice is that the observed t-value of the paired t-test and the t-test printed out by the linear model don't match:

```
t.test(rawRT ~ condition, bysubj, paired = TRUE)$statistic
```

```
##      t
## 3.109
```

```
round(summary(m0)$coefficients, 2)[, c(1:3)]
```

```
##              Estimate Std. Error t value
## (Intercept)      471.4       31.13  15.14
## conditionsubjgap -102.3       44.02  -2.32
```

This is because the linear model implements the unpaired (i.e., two sample) t-test:

```
summary(lm(rawRT ~ condition, bysubj))
```

```
##
## Call:
## lm(formula = rawRT ~ condition, data = bysubj)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -303.4  -116.4   -51.6    49.1   853.3
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)      471.4       31.1   15.14  <2e-16
## conditionsubjgap -102.3       44.0    -2.32   0.023
##
## Residual standard error: 202 on 82 degrees of freedom
## Multiple R-squared:  0.0618, Adjusted R-squared:  0.0503
## F-statistic:  5.4 on 1 and 82 DF,  p-value: 0.0226
```

```
round(t.test(rawRT ~ condition, bysubj, paired = FALSE)$statistic,
      2)
```

```
##      t
## 2.32
```

The paired t-test has an equivalent in the linear modeling framework: the linear mixed model. We turn next to this extension of the simple linear model. The command corresponding to the paired t-test in the linear modeling framework is:

```
m0.lmer <- lmer(rawRT ~ condition + (1 | subject),
  bysubj)
summary(m0.lmer)$coefficients
```

```
##              Estimate Std. Error t value
## (Intercept)      471.4       31.13  15.143
## conditionsubjgap -102.3       32.90  -3.109
```

To understand the connection between the paired t-test and the above command, it is necessary to consider how a paired t-test is assembled.

First, some background. If you have two random variables that have correlation ρ , the variance of the difference between the two random variables is:

$$\text{Var}(X_1 - X_2) = \text{Var}(X_1) + \text{Var}(X_2) - 2 \times \text{Cov}(X_1, X_2) \quad (3.18)$$

$\text{Cov}(X_1, X_2)$ is the covariance between the two random variables and is defined as:

$$\text{Cov}(X_1, X_2) = \rho \sqrt{\text{Var}(X_1)} \sqrt{\text{Var}(X_2)} \quad (3.19)$$

You can find the proofs of the above assertions in books like [Rice \(1995\)](#).

As discussed earlier, a paired t-test is used when you have paired data from subject $i = 1, \dots, n$ in two conditions, say conditions 1 and 2. Let's write the data as two vectors X_1, X_2 . Because the pairs of data points are coming from the same subject, they are

correlated with some correlation ρ . Assume that both conditions 1 and 2 have standard deviation σ .

To make this discussion concrete, let's generate some simulated bivariate data that are correlated. Assume that $\sigma = 1$, $\rho = 0.5$, and that the data are balanced.

```
library(MASS)
samplesize <- 12
mu <- c(0.3, 0.2)
rho <- 0.5
stddev <- 1
Sigma <- matrix(stddev, nrow = 2, ncol = 2) + diag(2)
Sigma <- Sigma/2
Sigma
```

```
##      [,1] [,2]
## [1,]  1.0  0.5
## [2,]  0.5  1.0
```

```
## simulated data:
x <- mvrnorm(n = samplesize, mu = mu, Sigma = Sigma,
             empirical = TRUE)
head(x)
```

```
##      [,1] [,2]
## [1,] -0.4241 -1.1746
## [2,]  0.4285  0.1175
## [3,] -1.1445 -1.6417
## [4,]  2.1234  0.1617
## [5,] -0.9182  0.3638
## [6,]  0.6915  1.6569
```

```
n <- samplesize
x1 <- x[, 1]
```

```
x2 <- x[, 2]
x1
```

```
## [1] -0.4241  0.4285 -1.1445  2.1234 -0.9182  0.6915
## [7] -0.3609  0.2910  0.9116  1.8084  0.4441 -0.2508
```

```
x2
```

```
## [1] -1.17458  0.11751 -1.64174  0.16166  0.36377
## [6]  1.65688  1.12933 -0.26363  1.21378  1.22185
## [11] -0.41870  0.03387
```

To carry out the paired t-test, we need to know the variance of $X_1 - X_2$ because the t-statistic will be:

$$t_{n-1} = \frac{X_1 - X_2}{\sqrt{\text{Var}(X_1 - X_2)/n}} \quad (3.20)$$

Now,

$$\text{Var}(X_1 - X_2) = \sigma^2 + \sigma^2 - 2\rho\sigma\sigma = 2\sigma^2(1 - \rho) \quad (3.21)$$

Now let's compute the t-statistic using the above formula. Let the actual data vectors be x_1, x_2 .

$$t_{n-1} = \frac{\text{mean}(x_1) - \text{mean}(x_2)}{\sqrt{\text{Var}(X_1 - X_2)/n}} \quad (3.22)$$

This simplifies to:

$$t_{n-1} = \frac{\text{mean}(x_1) - \text{mean}(x_2)}{\sqrt{2\sigma^2(1 - \rho)/n}} \quad (3.23)$$

Now compare the paired t-test output and the by-hand calculation:

```
t.test(x1, x2, paired = TRUE)$statistic
```

```
##      t
## 0.3464
```

```
(mean(x1) - mean(x2))/sqrt((2 * stddev^2 * (1 - rho))/n)
```

```
## [1] 0.3464
```

The linear mixed model we present next will fit exactly the same model as in the paired t-test above. To see this, suppose we have i subjects and $j = 1, 2$ conditions. For simplicity, assume that each subject sees each condition once (e.g., the by-subjects aggregated English relative clause data), so we have two data points from each subject. In other words, the data are paired.

Then, for condition 1, the dependent variable can be described by the equation:

$$y_{i1} = \beta_0 + u_{0i} + \varepsilon_{i1}$$

Here, β_0 is the mean reading time, and ε is the usual residual error term. The interesting new term is u_{0i} . This is the adjustment to the mean reading time for subject i . That is, if some subject is slower than average, u_{0i} will be a positive number; if a subject is faster than average, then that subject's adjustment u_{0i} will be negative in sign; and if a subject has exactly the same reading time as the mean for all subjects, then u_{0i} for that subject will be 0.

Similarly, for condition 2, the dependent variable is described by the equation:

$$y_{i2} = \beta_0 + \delta + u_{0i} + \varepsilon_{i2}$$

Here, δ is the additional time taken to process condition 2 (thus, this is the treatment contrast coding we saw earlier in this chapter).

If we subtract the equation for condition 2 from the equation for condition 1, the resulting equation is:

$$d_i = y_{i1} - y_{i2} = \delta + (\varepsilon_{i1} - \varepsilon_{i2})$$

The expectation of d_i is δ because the expectation of the ε terms is 0 (we set up the model such that $\varepsilon \sim \text{Normal}(0, \sigma)$).

Now, assuming that the error terms are correlated with correlation ρ , the result presented at the beginning of this section applies:

$$\text{Var}(y_{i1} - y_{i2}) = \sigma^2 + \sigma^2 - 2\rho\sigma^2 = 2\sigma^2(1 - \rho) \quad (3.24)$$

The generative distribution for d_i , the pairwise differences in the two conditions, is

$$d \sim \text{Normal}(\delta, \sqrt{2\sigma^2(1 - \rho)}) \quad (3.25)$$

But that is exactly the same standard deviation as the one used in the paired t-test.

So, the paired t-test will deliver exactly the same t-score as the above linear mixed model.

Let's check that this is true using our simulated data. In the code below, the term `(1|subj)` is the adjustment by subject to the intercepts—the term u_{0i} above.

```
library(lme4)
dat <- data.frame(y = c(x1, x2), cond = rep(letters[1:2],
  each = n), subj = rep(1:n, 2))
dat$cond <- factor(dat$cond)
contrasts(dat$cond) <- contr.sum(2)
contrasts(dat$cond)
```

```
##      [,1]
## a      1
## b     -1
```

```
summary(m <- lmer(y ~ cond + (1 | subj), dat))
```

```
## Linear mixed model fit by REML ['lmerMod']
## Formula: y ~ cond + (1 | subj)
##   Data: dat
##
## REML criterion at convergence: 65.6
##
## Scaled residuals:
##   Min      1Q  Median      3Q      Max
## -1.226 -0.533  0.032  0.487  1.737
##
## Random effects:
##   Groups      Name      Variance Std.Dev.
##   subj      (Intercept)  0.5      0.707
##   Residual                0.5      0.707
## Number of obs: 24, groups:  subj, 12
##
## Fixed effects:
##              Estimate Std. Error t value
## (Intercept)    0.250     0.250     1.00
## cond1          0.050     0.144     0.35
##
## Correlation of Fixed Effects:
##      (Intr)
## cond1 0.000
```

The t-statistic from the linear mixed model is exactly the same as that from the paired t-test.

With this as background, we are ready to look at linear mixed models in detail.

3.5 Linear mixed models

We return to our subject and object relative clause data from English (Grodner and Gibson, Expt 1). First we load the data as usual, define relative clause type as a sum coded predictor, and

create a new column called `so` that represents the contrast coding (± 1 sum contrasts), and a column that holds log-transformed reading time.

```
gg05e1 <- read.table("data/grodnergibsonE1crit.txt",
  header = TRUE)

gg05e1$so <- ifelse(gg05e1$condition == "objgap", 1,
  -1)
gg05e1$logrt <- log(gg05e1$rawRT)
```

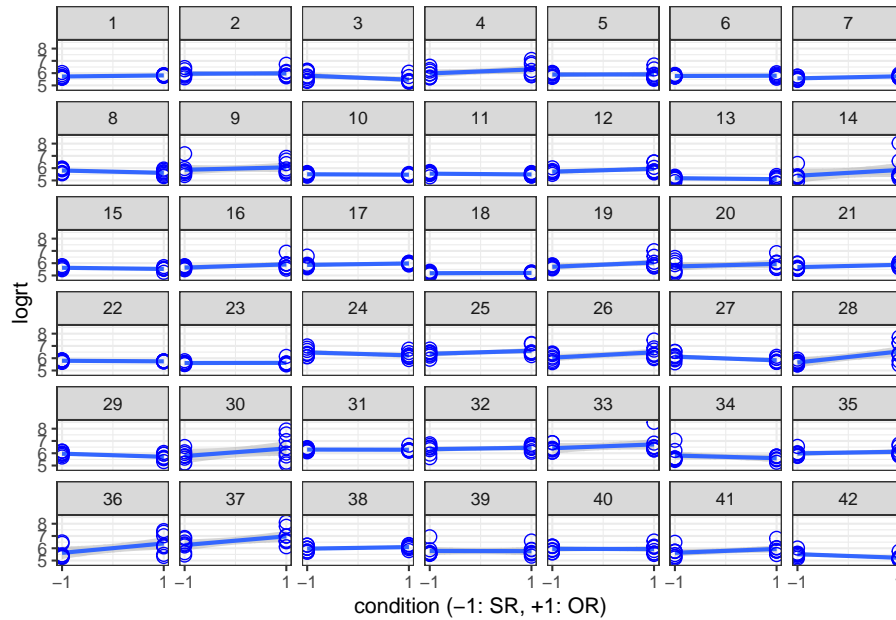
Recall that these data have multiple measurements from each subject for each condition:

```
t(xtabs(~subject + condition, gg05e1))

##           subject
## condition 1  2  3  4  5  6  7  8  9 10 11 12 13 14 15 16 17 18
##   objgap   8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8
##   subjgap   8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8
##           subject
## condition 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33
##   objgap    8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8
##   subjgap    8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8  8
##           subject
## condition 34 35 36 37 38 39 40 41 42
##   objgap    8  8  8  8  8  8  8  8  8
##   subjgap    8  8  8  8  8  8  8  8  8
```

We can visualize the different responses of subjects:

```
## `geom_smooth()` using formula 'y ~ x'
```



It's clear that different subjects have different effects of the relative clause manipulation: some slopes are positive sloping, some are flat, and some are negatively sloping. There is between-subject variability in the relative clause effect.

Given these differences between subjects, you could fit a separate linear model for each subject, collect together the intercepts and slopes for each subject, and then check if the slopes are significantly different from zero. There is a function in the package `lme4` that computes separate linear models for each subject: `lmList`.

```
library(lme4)
```

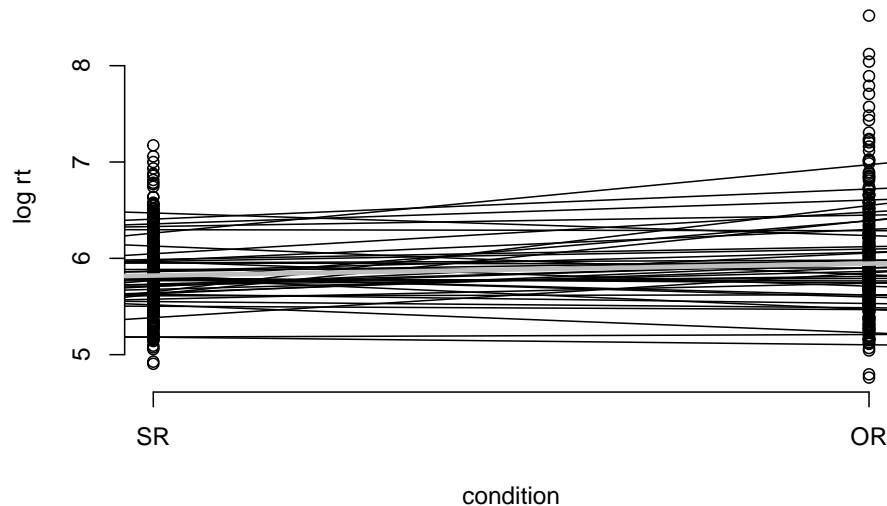
```
lm1list.fm1 <- lmList(logrt ~ so | subject, gg05e1)
```

One can extract the intercept and slope estimates for each subject. For example, for subject 1:

```
lm1list.fm1$`1`$coefficients
```

```
## (Intercept)          so
##      5.76962      0.04352
```

One can plot the individual lines for each subject, as well as the fit of a simple linear model `m0` for all the data taken together; this will show how each subject deviates in intercept and slope from the model `m0`'s intercept and slope.



To find out if there is an effect of relative clause type, we simply need to check whether the slopes of the individual subjects' fitted lines taken together are significantly different from zero. A one-sample t-test will achieve this:

```
t.test(coef(lm1)[2])
```

```
##
## One Sample t-test
##
## data:  coef(lm1)[2]
## t = 2.8, df = 41, p-value = 0.008
## alternative hypothesis: true mean is not equal to 0
## 95 percent confidence interval:
##  0.01745 0.10658
## sample estimates:
```

```
## mean of x
## 0.06202
```

The above test is *exactly* the same as the paired t-test and the varying intercepts linear mixed model that we fit in the last chapter using the by-subject aggregated data:

```
bysubj <- aggregate(log(rawRT) ~ subject + condition,
  mean, data = gg05e1)

colnames(bysubj)[3] <- "logrt"

t.test(logrt ~ condition, bysubj, paired = TRUE)$statistic
```

```
## t
## 2.81
```

```
## compare with linear mixed model:
summary(lmer(logrt ~ condition + (1 | subject), bysubj))$coefficients[2,
]
```

```
## Estimate Std. Error t value
## -0.12403 0.04414 -2.81021
```

The above `lmList` model we just fit is called *repeated measures regression*. We now look at how to model unaggregated data using the linear mixed model. Incidentally, this repeated measures regression model is now only of historical interest, and useful only for understanding the linear mixed model, which is the modern standard approach.

We turn next to three main types of linear mixed model; other variants will be introduced in later chapters.

3.5.1 Model type 1: Varying intercepts

The *linear mixed model* does something related to the above by-subject fits, but with some crucial twists, as we see below. In the model shown below, the statement

$$(1 \mid \text{subject}) \quad (3.26)$$

adjusts the grand mean estimates of the intercept by a term (a number) for each subject.

```
m0.lmer <- lmer(logrt ~ so + (1 | subject), gg05e1)
```

Notice that we did not aggregate the data.

Here is the abbreviated output:

Random effects:

Groups	Name	Variance	Std.Dev.
subject	(Intercept)	0.09983	0.3160
Residual		0.14618	0.3823

Number of obs: 672, groups: subject, 42

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	5.88306	0.05094	115.497
so	0.06202	0.01475	4.205

One thing to notice is that the coefficients (intercept and slope) of the fixed effects of the above model are identical to those in the linear model `m0` above. What is different between the linear model and the linear mixed model is the standard error. In the latter, the standard error is determined by more than one source of variance, as we explain below.

The intercept adjustments for each subject can be viewed by typing:

```
## first 10 subjects' intercept adjustments:
```

```
ranef(m0.lmer)$subject[, 1][1:10]
```

```
## [1] -0.103928  0.077195 -0.230621  0.234198  0.008828
```

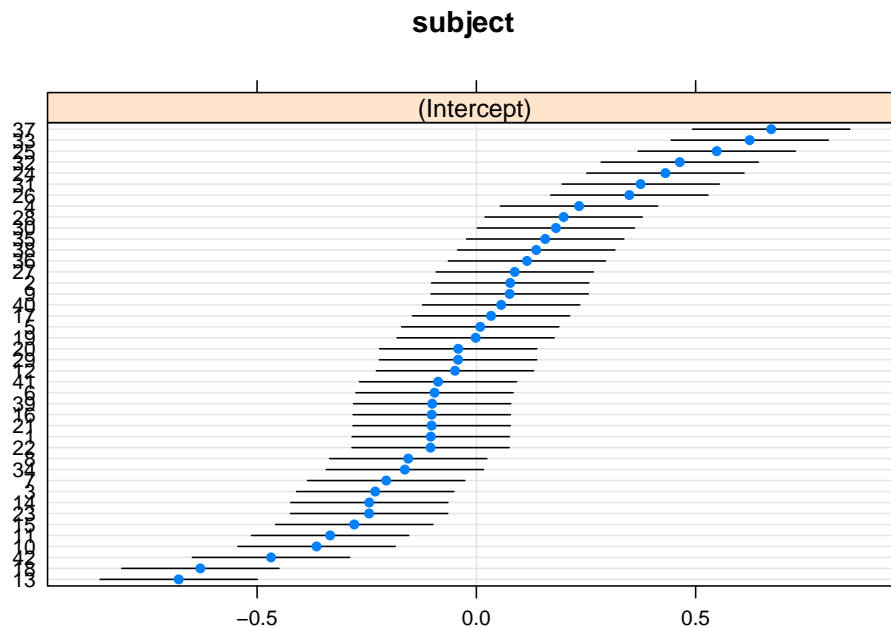
```
## [6] -0.095363 -0.205571 -0.155371  0.075944 -0.364367
```

Here is another way to summarize the adjustments to the grand mean intercept by subject. The error bars represent 95% confidence intervals.

```
library(lattice)
```

```
print(dotplot(ranef(m0.lmer, condVar = TRUE)))
```

```
## $subject
```



3.5.2 The formal statement of the varying intercepts model

The model `m0.lmer` above prints out the following type of linear model. i indexes subject, and j indexes items.

Once we know the subject id and the item id, we know which subject saw which condition:

```
subset(gg05e1, subject == 1 & item == 1)
```

```
##   subject item condition rawRT so logrt
## 6         1     1    objgap   320   1 5.768
```

The mathematical form of the linear mixed model is:

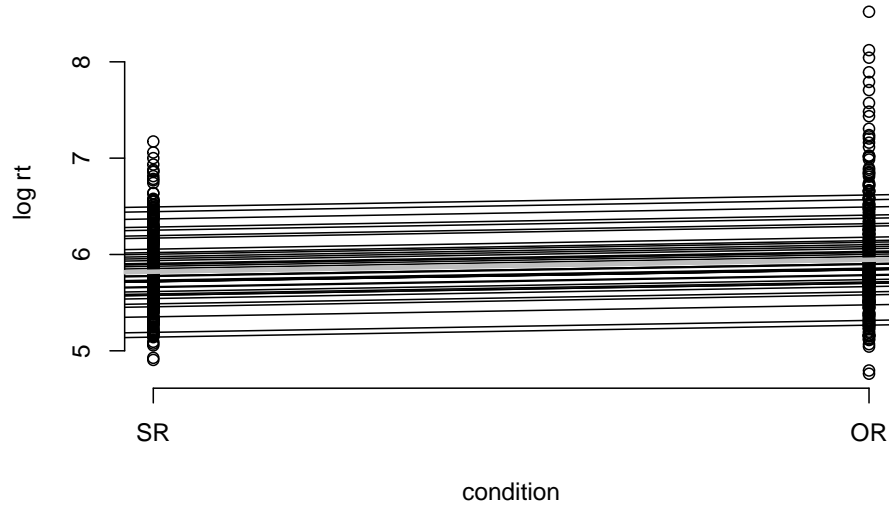
$$y_{ij} = \beta_0 + u_{0i} + \beta_1 \times so_{ij} + \varepsilon_{ij} \quad (3.27)$$

The *only* new thing here beyond the linear model we saw earlier is the by-subject adjustment to the intercept. These by-subject adjustments to the intercept u_{0i} are assumed by `lmer` to come from a normal distribution centered around 0:

$$u_{0i} \sim Normal(0, \sigma_{u0}) \quad (3.28)$$

The ordinary linear model `m0` has one intercept β_0 for all subjects, whereas this linear mixed model with varying intercepts `m0.lmer` has a different intercept ($\beta_0 + u_{0i}$) for each subject i .

We can visualize the adjustments for each subject to the intercepts as shown below.



An important point is that in this model there are two variance components or sources of variance (cf. the linear model, which had only one):

- $u_0 \sim \text{Normal}(0, \sigma_{u0})$
- $\varepsilon \sim \text{Normal}(0, \sigma)$

These two standard deviations determine the standard error of the β_1 slope parameter.

3.5.3 Model type 2: Varying intercepts and slopes, without a correlation

Unlike the figure associated with the `lmlist.fm1` model above, which also involves fitting separate models for each subject, the model `m0.lmer` assumes *different intercepts* for each subject *but the same slope*.

We can choose to fit different intercepts as well as different slopes for each subject. To achieve this, assume now that each subject's slope is also adjusted by subject:

$$y_{ij} = \beta_0 + u_{0i} + (\beta_1 + u_{1i}) \times so_{ij} + \varepsilon_{ij} \quad (3.29)$$

That is, we additionally assume that $u_{1i} \sim \text{Normal}(0, \sigma_{u1})$. The `lmer` notation for fitting separate intercepts and slopes is

(1+so||subject). We will just explain what the double vertical bars represent.

```
m1.lmer <- lmer(logrt ~ so + (1 + so || subject), gg05e1)
```

The output of this model will now show that there are not two but three sources of variability. These are:

- $u_0 \sim \text{Normal}(0, \sigma_{u0})$
- $u_1 \sim \text{Normal}(0, \sigma_{u1})$
- $\varepsilon \sim \text{Normal}(0, \sigma)$

In particular, the model estimates the following standard deviations:

- $\hat{\sigma}_{u0} = 0.317$
- $\hat{\sigma}_{u1} = 0.110$
- $\hat{\sigma} = 0.365$.

Random effects:

Groups	Name	Variance	Std.Dev.
subject	(Intercept)	0.1006	0.317
subject.1	so	0.0121	0.110
Residual		0.1336	0.365

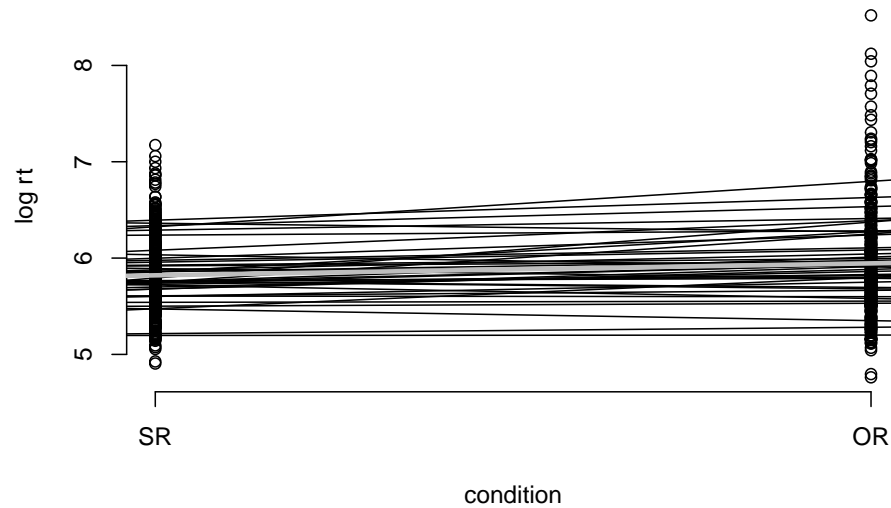
Number of obs: 672, groups: subject, 42

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	5.8831	0.0509	115.50
so	0.0620	0.0221	2.81

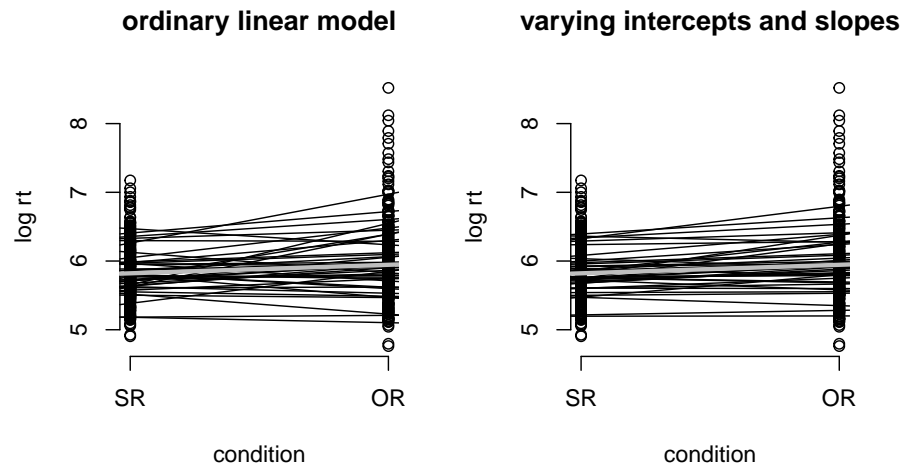
These fits for each subject are visualized below (the gray line shows the model with a single intercept and slope, i.e., our old model `m0`):

varying intercepts and slopes for each subject



3.5.3.1 Comparing `lmList` model with the varying intercepts model

Compare this model with the `lmList.fm1` model we fitted earlier:



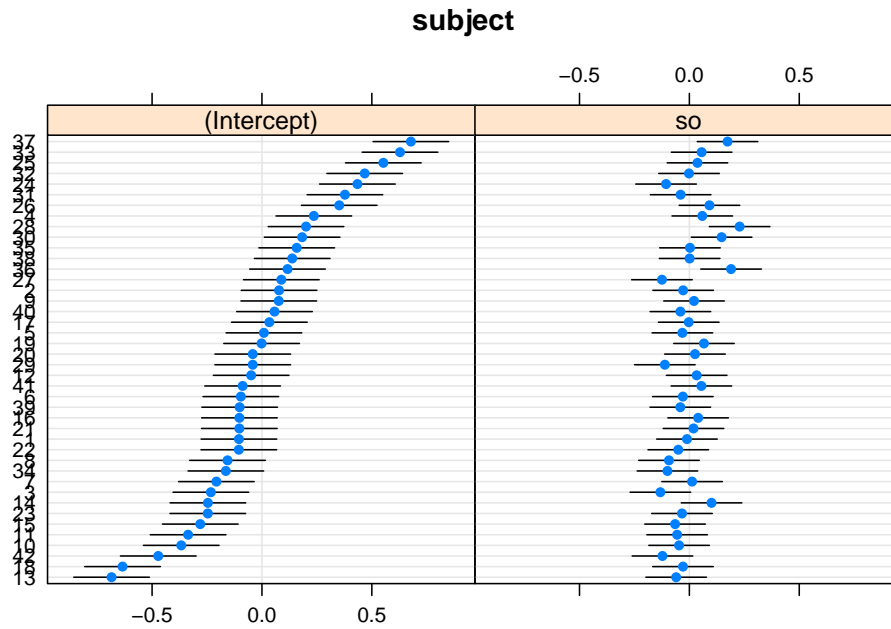
What is striking is that each subject's estimated best fit line is "smoothened out" compared to the `lmList` fits. This aspect of the linear mixed model is called shrinkage; we return to this point presently.

3.5.3.2 Visualizing random effects

As before, it is instructive to visualize the individual level adjustments to the intercept and slope:

```
print(dotplot(ranef(m1.lmer, condVar = TRUE)))
```

```
## $subject
```



What this is showing is wide variability in the mean reading times between subjects, but very little variation in the slope between subjects.

3.5.3.3 The formal statement of varying intercepts and varying slopes linear mixed model

Here is the full statement of the varying intercept and slopes model. Again, i indexes subjects, j items.

$$y_{ij} = \beta_0 + u_{0i} + (\beta_1 + u_{1i}) \times so_{ij} + \varepsilon_{ij} \quad (3.30)$$

As mentioned before, there are now three variance components:

- $u_0 \sim \text{Normal}(0, \sigma_{u0})$
- $u_1 \sim \text{Normal}(0, \sigma_{u1})$
- $\varepsilon \sim \text{Normal}(0, \sigma)$

3.5.3.4 Crossed random effects for subjects and for items

Now, one interesting fact about the varying intercepts and slopes model is that it doesn't capture all the sources of variance yet. The items also contribute sources of variance: just like subjects, items may also have different intercepts and slopes.

Notice that subjects and items are fully crossed: each subject sees each item once.

```
head(xtabs(~subject + item, gg05e1))
```

```
##          item
## subject 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16
##      1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      3 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      4 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      5 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      6 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
```

Linear mixed model with crossed subject and items random effects can be defined with the following syntax:

```
m2.lmer <- lmer(logrt ~ so + (1 + so || subject) +
  (1 + so || item), gg05e1)
```

Analogously to the preceding example, now there are five variance components:

Random effects:

Groups	Name	Variance	Std.Dev.
subject	(Intercept)	0.10090	0.3177
subject.1	so	0.01224	0.1106
item	(Intercept)	0.00127	0.0356

```

item.1    so          0.00162  0.0402
Residual              0.13063  0.3614
Number of obs: 672, groups:  subject, 42; item, 16

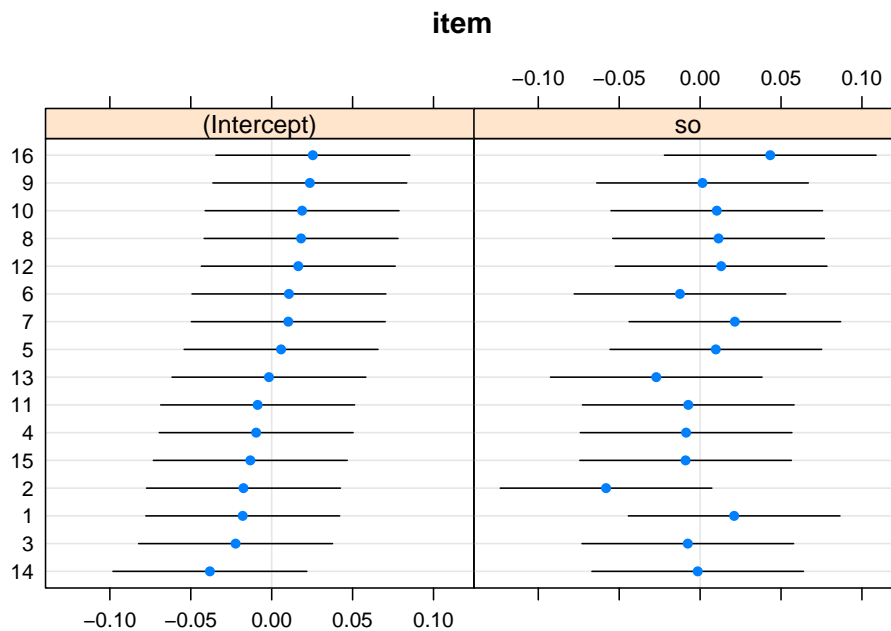
```

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	5.8831	0.0517	113.72
so	0.0620	0.0242	2.56

The item intercept and slope adjustments can be visualized as well. Notice that there is a lot less item-level variation; this is often the case in planned experiments like this one, where the experimental items are carefully constructed to vary as little as possible.

```
print(dotplot(ranef(m2.lmer, condVar = TRUE))$item)
```



One thing missing in the above models is any assumption about the relationship between the intercept and slope adjustments by subject and by item. It is possible that the intercept and slope adjustments are correlated: e.g., there could be a theoretical prediction that says that the slower a subject's average reading time,

the larger the difference between object and subject relative clause reading times. This kind of prediction can be tested by testing what the correlation is between the varying intercepts and slopes. We turn to this model next.

3.5.4 Model type 3: Varying intercepts and varying slopes, with correlation

A correlation can be introduced between the intercept and slope adjustments by using a single vertical bar instead of two vertical bars in the random effects structure:

```
m3.lmer <- lmer(logrt ~ so + (1 + so | subject) + (1 +
  so | item), gg05e1)
```

```
## boundary (singular) fit: see ?isSingular
```

To understand what this model is doing, we have to recall what a bivariate/multivariate distribution is.

Random effects:

Groups	Name	Variance	Std.Dev.	Corr
subject	(Intercept)	0.10103	0.3178	
	so	0.01228	0.1108	0.58
item	(Intercept)	0.00172	0.0415	
	so	0.00196	0.0443	1.00 <= degeneracy
Residual		0.12984	0.3603	

Number of obs: 672, groups: subject, 42; item, 16

Fixed effects:

	Estimate	Std. Error	t value
(Intercept)	5.8831	0.0520	113.09
so	0.0620	0.0247	2.51

The correlations (0.58 and 1.00) you see in the model output below are the correlations between the varying intercepts and slopes for subjects and for items. Notice that the variance covariance matrix for items is degenerate: its correlation is 1. This matrix cannot be inverted.

When the correlation is +1 or -1 or near these numbers, this means that the optimizer in lme4 is unable to estimate the correlation parameter, usually due to there not being enough data. If you are in such a situation, you are better off not trying to estimate this parameter with the data you have, and instead fitting one of the simpler models. We will return to this point when discussing model selection. For further discussion, see [Barr et al. \(2013\)](#), [Bates et al. \(2015\)](#), and [Matuschek et al. \(2017\)](#).

3.5.4.1 Formal statement of varying intercepts and varying slopes linear mixed model with correlation

As usual, i indexes subjects, j items. The vector `so` is the sum-coded factor levels: +1 for object relatives and -1 for subject relatives. The only new thing in this model is the item-level effects, and the specification of the variance-covariance matrix for subjects and items, in order to include the correlation parameters.

$$y_{ij} = \alpha + u_{0i} + w_{0j} + (\beta + u_{1i} + w_{1j}) \times so_{ij} + \varepsilon_{ij} \quad (3.31)$$

where $\varepsilon_{ij} \sim \text{Normal}(0, \sigma)$ and

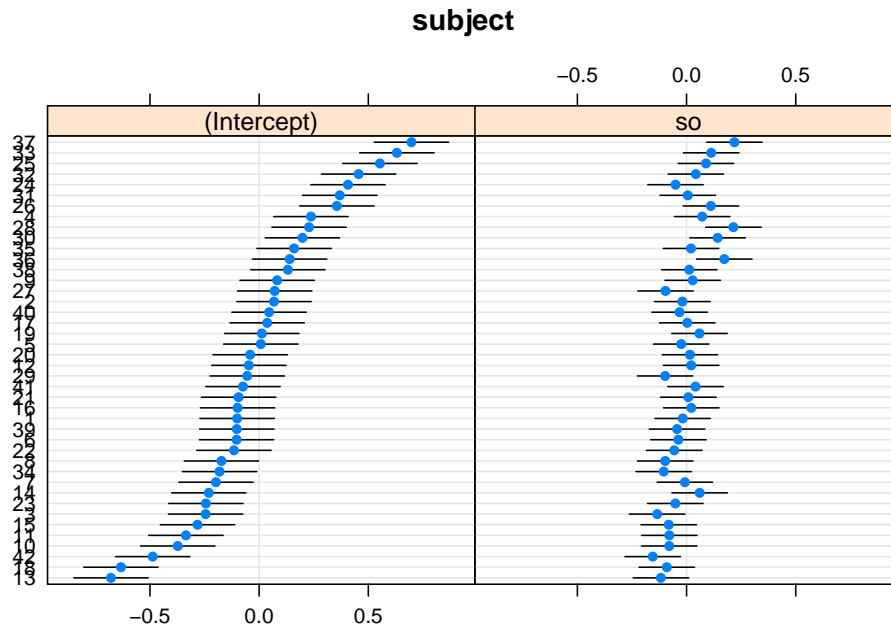
$$\Sigma_u = \begin{pmatrix} \sigma_{u0}^2 & \rho_u \sigma_{u0} \sigma_{u1} \\ \rho_u \sigma_{u0} \sigma_{u1} & \sigma_{u1}^2 \end{pmatrix} \quad \Sigma_w = \begin{pmatrix} \sigma_{w0}^2 & \rho_w \sigma_{w0} \sigma_{w1} \\ \rho_w \sigma_{w0} \sigma_{w1} & \sigma_{w1}^2 \end{pmatrix} \quad (3.32)$$

$$\begin{pmatrix} u_0 \\ u_1 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma_u \right), \quad \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma_w \right) \quad (3.33)$$

3.5.4.2 Visualizing the random effects

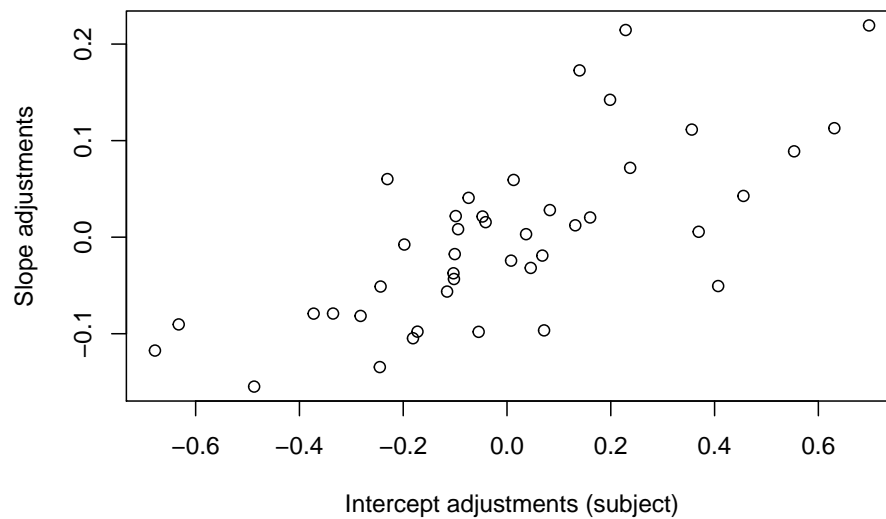
One can visualize the correlation between intercepts and slopes by subjects. The positive correlation of 0.58 between subject intercept and slope adjustments implies that slower subjects show larger effects. However, the dotplot below doesn't show a convincing pattern:

```
print(dotplot(ranef(m3.lmer, condVar = TRUE))$subject)
```



The correlation pattern is easier to see if we plot the slope adjustments against the intercept adjustments.

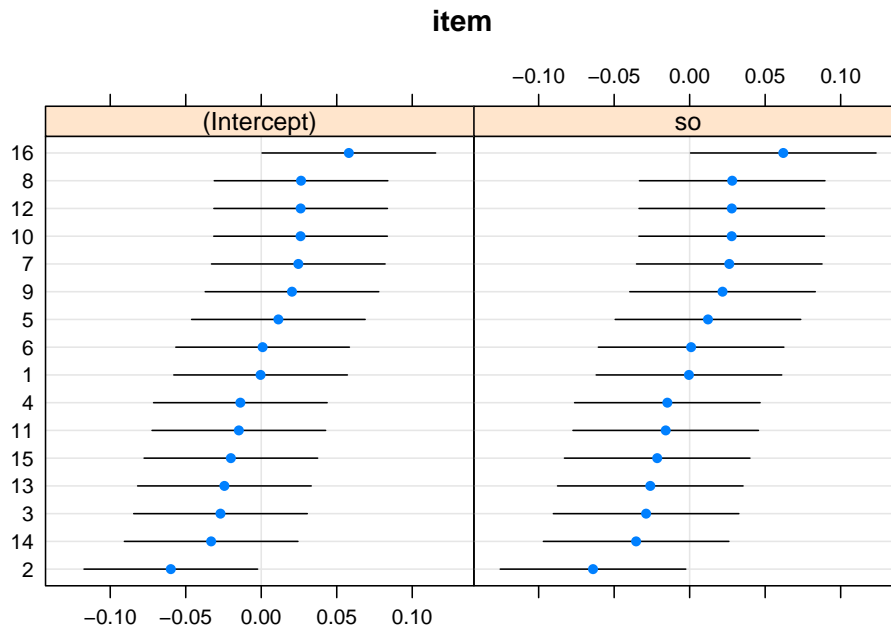
```
plot(ranef(m3.lmer)$subject[, 1], ranef(m3.lmer)$subject[,
  2], xlab = "Intercept adjustments (subject)", ylab = "Slope adjustments")
```



When we talk about hypothesis testing, we will look at what inferences we can draw from this correlation.

The dotplot showing the item-level effects shows a perfect correlation between intercept and slope adjustments, but as mentioned above these are from a degenerate variance covariance matrix and not meaningful.

```
print(dotplot(ranef(m3.lmer, condVar = TRUE))$item)
```



3.6 Shrinkage in linear mixed models

The estimate of the effect by participant is smaller than when we fit a separate linear model to the subject's data. This is called shrinkage in linear mixed models: the individual level estimates are shrunk towards the mean slope. The less data we have from a given subject, the greater the shrinkage.

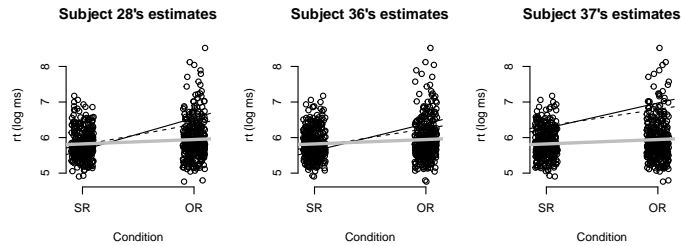


FIGURE 3.1: The figures show linear model fits (the grand mean estimates) for three subjects; shown are the simple linear model fit on all the data (gray line), the `lmList` model fit to the individual subject's data (black line), and the linear mixed model fit (the broken line). In all three subjects' models, the linear mixed model estimates are shrunk towards the grand mean (gray line) estimates.

3.6.0.1 Shrinkage in action: when data are missing

The importance and value of shrinkage becomes clear once we simulate a situation where there is some missing data. Missingness can happen in experiments, either due to lost measurements (arising from computer error or programming errors), or some other reason. To see what happens when we have missing data, let's randomly delete some data from one subject. We will randomly delete 50% of subject 37's data:

```
set.seed(4321)
## choose some data randomly to remove:
rand <- rbinom(1, n = 16, prob = 0.5)
```

Here are subject 37's reading times (16 data points):

```
gg05e1[which(gg05e1$subject == 37), ]$rawRT
```

```
## [1] 770 536 686 578 457 487 2419 884 3365 233
## [11] 715 671 1104 281 1081 971
```

Now, we randomly delete half the data:

```
gg05e1$deletedRT <- gg05e1$rawRT
gg05e1[which(gg05e1$subject == 37), ]$deletedRT <- ifelse(rand,
  NA, gg05e1[which(gg05e1$subject == 37), ]$rawRT)
```

Now subject 37's estimates are going to be pretty wild, because they are based on much less data (even one extreme value can strongly influence the mean):

```
subset(gg05e1, subject == 37)$deletedRT
```

```
## [1] 770 NA 686 578 NA NA NA NA 3365 233
## [11] NA 671 1104 NA NA 971
## [1] 6.617
## [1] 0.3554
## [1] 6.688
## [1] 0.3884
```

Now fit the hierarchical model and examine subject 37's estimates on undeleted vs deleted data:

What we see here is that the estimates from the hierarchical model are barely affected by the missingness, but the estimates from the `lmList` model are heavily affected. This means that linear mixed models will give you more robust estimates (think Type M error!) compared to no pooling models. This property of shrinkage is one reason why linear mixed models are so important in cognitive science.

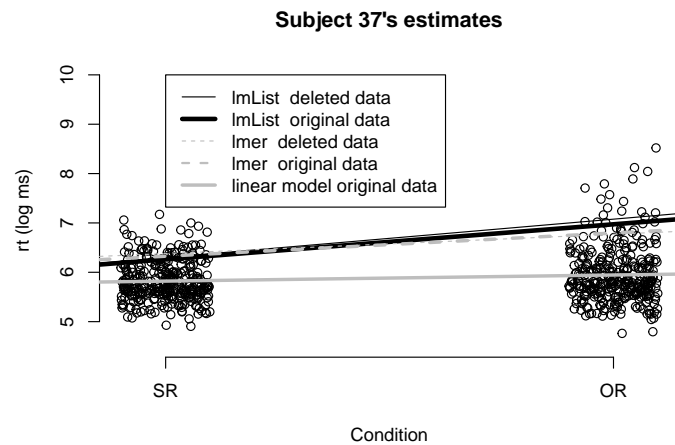


FIGURE 3.2: The figures show linear model fits (the grand mean estimates) for subject 37. When using `lmList`, deleting data leads to very different estimates; but using `lmer`, deleting half the data from this subject hardly affects the individual subject's estimates.

3.7 Summary

3.8 Exercises

Download the data-set `E1_data.csv`. Then run the following commands to load the `lme4` library and to set up your data for analysis:

```
library(lme4)

## load data:
dat <- read.csv("data/E1_data.csv", header = TRUE)
## convert RT to milliseconds:
dat$RT <- dat$RT * 1000
## choose critical region:
word_n <- 4
```

```
## subset critical data:  
crit <- subset(dat, Position == word_n)
```

The data consist of a repeated measures experiment comparing two conditions which are labeled Type 1 and Type 2. The column Sub refers to subject id, and the column ID refers to item id. RT refers to reading time in seconds (we have converted it above to milliseconds); NA is missing data. You can ignore the other columns. This is a standard Latin square design. We will work with the data frame `crit` below.

3.8.1 By-subjects t-test

Using RT as a dependent variable, carry out the appropriate by-subjects t-test to evaluate the null hypothesis that there is no difference between the conditions labeled Type 1 and 2. Write down all the R commands needed to do the appropriate t-test, and the resulting t-value and p-value. State whether we can reject the null hypothesis given the results of the t-test; explain why.

3.8.2 Fitting a linear mixed model

Now, using the data-frame called `crit` above as the data-frame, fit a linear mixed model (called M0) with the column called Type coded as sum contrasts.

Assume varying intercepts for subjects and varying intercepts for items (varying intercepts are sometimes called random intercepts). Write down the linear mixed models command, and write down the fixed-effects estimates (Intercept and slope) along with their standard errors. State whether we can reject the null hypothesis given the results of the t-value shown in the linear mixed model output; explain why.

3.8.3 t-test vs. linear mixed model

Why do the results of the t-test and the linear mixed model M0 differ?

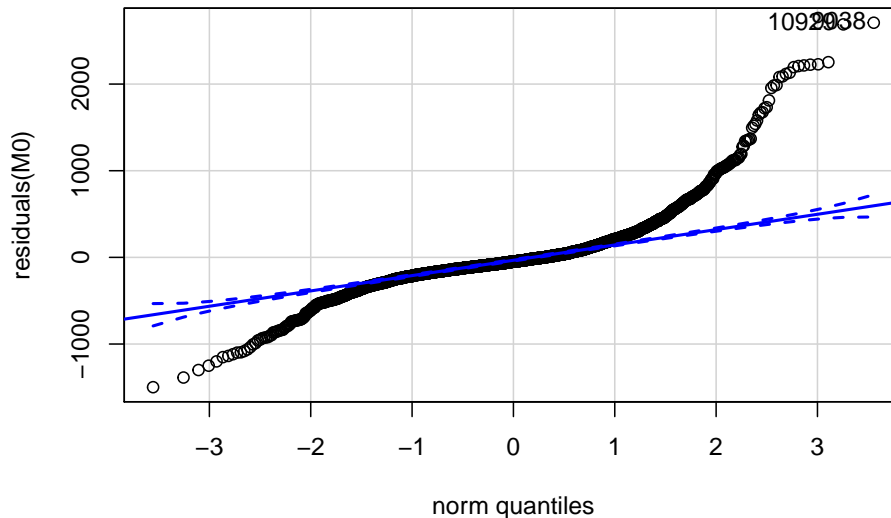
3.8.4 Power calculation using `power.t.test`

The researcher wants to achieve 80% statistical power in a future study. Based on the available data above, she determines that the standard error (note: not the standard deviation!) of the difference in means between the conditions Type 1 and Type 2 is 21. She has reason to believe that the true difference in means is 30 ms. What is the number of participants (to the nearest whole number) needed to achieve approximately 80% power? Use the `power.t.test` function to compute your answer. Write down the `power.t.test` function specification you used, as well as the number of participants needed, based on the output of the `power.t.test` function.

3.8.5 Residuals

The plot below shows the distribution of the residuals from model M0 plotted against the standard normal distribution with mean 0 and standard deviation 1. Explain what the plot tells us about one of the model assumptions of the linear mixed model M0 that you fit earlier.

(You can ignore the numbers below the plot.)



9038 10929 726 2544

3.8.6 Understanding contrast coding

Using only your estimates of the intercept and the slope in model M0's fixed effects output, write down the mean of the condition labeled Type 1 in the data, and the mean of the condition labeled Type 2.

3.8.7 Understanding the fixed-effects output

Suppose that the model M0's output for the fixed effects analysis were as follows. SO is a sum-coded contrast specification for the conditions labeled Type.

```
results
```

##	Estimate	Std. Error	t value
## (Intercept)	686.01	47.54	14.43
## SO	18.94	NA	2.00

What is the value of the standard error of the slope (SO), which is labeled NA above?

3.8.8 Understanding the null hypothesis test

A researcher fits a linear mixed model to compare the reading times between two conditions (a) and (b), just like in the above study. Her hypothesis is that the mean for condition (a) is larger than the mean for (b). She observes that condition a has sample mean 500 ms, and condition (b) has sample mean 450 ms. She also establishes from the linear mixed model that the t-value is 1.94. The approximate p-value associated with this t-value is 0.052. Answer the following: (A) Do we have evidence against the null hypothesis and (B) do we have evidence for the particular research hypothesis that the researcher has?

The researcher runs the same analysis as above on a new data-set that has the same design as above, and now gets a p-value of 0.001. Now she has stronger evidence than in the above case where the p-value was 0.052. What does she have stronger evidence for?



4

Hypothesis testing using the likelihood ratio test

We started the book with the one-sample t-test. There, we had the following procedure:

- Given independent and identically distributed data y , define a null hypothesis: $H_0 : \mu = \mu_0$
- Compute the sample mean \bar{y} and the standard error SE
- Reject the null hypothesis if the absolute value of \bar{y}/SE is larger than 2.

Here, we turn to a closely related test: the *likelihood ratio test statistic*.

4.1 The likelihood ratio test: The theory

Suppose that X_1, \dots, X_n are independent and normally distributed with mean μ and standard deviation σ (assume for simplicity that σ is known).

Let the null hypothesis be $H_0 : \mu = \mu_0$ and the alternative be $H_1 : \mu \neq \mu_0$. Here, μ_0 is a number, such as 0.

The likelihood of the data y can be computed under the null model, in which $\mu = \mu_0$, and under the alternative model, in which μ has some specific alternative value. To make this concrete, imagine 10 data points being generated from a Normal(0,1).

```
y <- rnorm(10)
```

We can compute the joint likelihood under a null hypothesis that $\mu = 0$:

```
likNULL <- prod(dnorm(y, mean = 0, sd = 1))
```

On the log scale, we would need to add the log likelihoods of each data point:

```
loglikNULL <- sum(dnorm(y, mean = 0, sd = 1, log = TRUE))  
loglikNULL
```

```
## [1] -11.6
```

Similarly, we can compute the log likelihood with μ equal to the maximum likelihood estimate of μ , the sample mean.

```
loglikALT <- sum(dnorm(y, mean = mean(y), sd = 1, log = TRUE))  
loglikALT
```

```
## [1] -11.59
```

Essentially, the likelihood ratio test compares the ratio of likelihoods of the two models; on the log scale, the difference in log likelihood is taken. The likelihood ratio test then chooses the model with the higher log likelihood, provided that the higher likelihood is high enough (we will just make this more precise).

One can specify the test in general terms as follows. Suppose that the likelihood is with respect to some parameter θ . We can evaluate the likelihood at μ_0 , the null hypothesis value of the parameter, and evaluate the likelihood using the maximum likelihood estimate $\hat{\theta}$ of the parameter. The likelihood ratio can then be written as follows:

$$\Lambda = \frac{\max_{\theta \in \omega_0} (lik(\theta))}{\max_{\theta \in \omega_1} (lik(\theta))} \quad (4.1)$$

where, $\omega_0 = \{\mu_0\}$ and $\omega_1 = \{\forall \mu \mid \mu \neq \mu_0\}$. The function max just selects the maximum value of any choices of parameter values; in the case of the null hypothesis there is only one value, μ_0 . In the case of the alternative model, the maximum likelihood estimate $\hat{\theta}$ is the maximum value.

Now, assuming that the data are coming from a normal distribution, the numerator of the likelihood ratio statistic is:

$$\frac{1}{(\sigma\sqrt{2\pi})^n} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu_0)^2 \right) \quad (4.2)$$

For the denominator, the MLE \bar{X} is taken as μ :

$$\frac{1}{(\sigma\sqrt{2\pi})^n} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \right) \quad (4.3)$$

The likelihood ratio statistic is then:

$$\Lambda = \frac{\frac{1}{(\sigma\sqrt{2\pi})^n} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu_0)^2 \right)}{\frac{1}{(\sigma\sqrt{2\pi})^n} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \right)} \quad (4.4)$$

Canceling out common terms:

$$\Lambda = \frac{\exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu_0)^2 \right)}{\exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \right)} \quad (4.5)$$

Taking logs:

$$\begin{aligned}
\log \Lambda &= \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu_0)^2 \right) - \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \bar{X})^2 \right) \\
&= -\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (X_i - \mu_0)^2 - \sum_{i=1}^n (X_i - \bar{X})^2 \right)
\end{aligned} \tag{4.6}$$

Now, it is a standard algebraic trick to rewrite $\sum_{i=1}^n (X_i - \mu_0)^2$ as a sum of two terms:

$$\sum_{i=1}^n (X_i - \mu_0)^2 = \sum_{i=1}^n (X_i - \bar{X})^2 + n(\bar{X} - \mu_0)^2 \tag{4.7}$$

If we rearrange terms, we obtain:

$$\sum_{i=1}^n (X_i - \mu_0)^2 - \sum_{i=1}^n (X_i - \bar{X})^2 = n(\bar{X} - \mu_0)^2 \tag{4.8}$$

Now, we just established above that $\log \Lambda$ is:

$$\log \Lambda = -\frac{1}{2\sigma^2} \left(\sum_{i=1}^n (X_i - \mu_0)^2 - \sum_{i=1}^n (X_i - \bar{X})^2 \right) \tag{4.9}$$

Consider the term in the brackets:

$$\left(\sum_{i=1}^n (X_i - \mu_0)^2 - \sum_{i=1}^n (X_i - \bar{X})^2 \right) \tag{4.10}$$

This can be rewritten as:

$$n(\bar{X} - \mu_0)^2 \tag{4.11}$$

Rewriting in this way gives us:

$$\ell = -\frac{1}{2\sigma^2} n(\bar{X} - \mu_0)^2 \tag{4.12}$$

Rearranging terms:

$$-2\ell = \frac{n(\bar{X} - \mu_0)^2}{\sigma^2} \quad (4.13)$$

Or even more transparently:

$$-2\ell = \frac{(\bar{X} - \mu_0)^2}{\frac{\sigma^2}{n}} \quad (4.14)$$

This should remind you of the t-test! Basically, just like in the t-test, what this is saying is that we reject the null when $|\bar{X} - \mu_0|$, or negative two times the difference in log likelihood, is large!

Now we will define what it means for -2ℓ to be large. We will define the *likelihood ratio test statistic* as follows. Here, $Lik(\theta)$ refers to the likelihood given some value θ for the parameter, and $\log Lik(\theta)$ refers to the log likelihood.

$$\begin{aligned} \Lambda &= -2 \times (Lik(\theta_0)/Lik(\theta_1)) \\ \log \Lambda &= -2 \times \{\log Lik(\theta_0) - \log Lik(\theta_1)\} \end{aligned} \quad (4.15)$$

where θ_1 and θ_0 are the estimates of θ under the alternative and null hypotheses, respectively. The likelihood ratio test rejects H_0 if $\log \Lambda$ is sufficiently large. As the sample size approaches infinity, $\log \Lambda$ approaches the chi-squared distribution:

$$\log \Lambda \rightarrow \chi_r^2 \text{ as } n \rightarrow \infty \quad (4.16)$$

where r is called the degrees of freedom and is the difference in the number of parameters under the null and alternative hypotheses.

The above result is called *Wilks' theorem*. The proof of Wilks' theorem is fairly involved but you can find it in Lehmann's textbook *Testing Statistical Hypotheses*.

Note that sometimes you will see the form:

$$\log \Lambda = 2\{\ell(\theta_1) - \ell(\theta_0)\} \quad (4.17)$$

It should be clear that both statements are saying the same thing; in the second case, we are just subtracting the null hypothesis log likelihood from the alternative hypothesis log likelihood, so the negative sign disappears.

That's the theory. Let's see how the likelihood ratio test works for (a) simulated data, and (b) our running example, the English relative clause data from [Grodner and Gibson \(2005\)](#).

4.2 A practical example using simulated data

A practical example will make the usage of this test clear. Let's just simulate a linear model:

```
x <- 1:10
y <- 10 + 20 * x + rnorm(10, sd = 10)
```

Here, the null hypothesis that the slope is 0 is false (it has value 20). Now, we fit a null hypothesis model, without a slope:

```
## null hypothesis model:
m0 <- lm(y ~ 1)
```

We will compare this model's log likelihood with that of the alternative model, which includes an estimate of the slope:

```
## alternative hypothesis model:
m1 <- lm(y ~ x)
```

The difference in log likelihood, multiplied with -2, is:

```
LogLambda <- -2 * (logLik(m0) - logLik(m1))
```



```
## observed value:  
LogLambda[1]
```

```
## [1] 34.49
```

The difference in the number of parameters in the two models is one, so $-2 \log \Lambda$ has the distribution χ_1^2 . Is the observed value of $-2 \log \Lambda$ unexpected under this distribution? We can calculate the probability of obtaining the likelihood ratio statistic we observed above, or a value more extreme, given the χ_1^2 distribution.

```
pchisq(LogLambda[1], df = 1, lower.tail = FALSE)
```

```
## [1] 4.286e-09
```

Just like the critical t-value in the t-test, the critical chi-squared value here is:

```
## critical value:  
qchisq(0.95, df = 1)
```

```
## [1] 3.841
```

If minus two times the observed difference in log likelihood is larger than this critical value, we reject the null hypothesis.

Note that in the likelihood test above, we are comparing one nested model against another: the null hypothesis model is nested inside the alternative hypothesis model. What this means is that the alternative hypothesis model contains all the parameters in the null hypothesis model (i.e., the intercept) plus another one (the slope).

4.3 A real-life example: The English relative clause data

The likelihood ratio test is also the way that hypothesis testing is done with the linear mixed model. Here is how it works. Let's look again at the [Grodner and Gibson \(2005\)](#) English relative clause data. The null hypothesis here refers to the slope parameter. When we have the sum contrast coding, the intercept β_0 refers to the grand mean, and the slope β_1 is the amount by which subject and object relative clause mean reading times deviate from the grand mean. Testing the null hypothesis that β_0 is 0 amounts to testing whether there is any difference in means between the two relative clause types. This becomes clear if we consider the following.

Let object relatives be coded as +1 and subject relatives as -1. Then, the mean reading time μ_{or} for object relatives in the linear mixed model is:

$$\mu_{or} = \beta_0 + \beta_1 \quad (4.18)$$

Similarly, the mean reading time μ_{sr} for subject relatives is:

$$\mu_{sr} = \beta_0 - \beta_1 \quad (4.19)$$

If the null hypothesis is that $\mu_{or} - \mu_{sr} = 0$, then this amounts to saying that:

$$(\beta_0 + \beta_1) - (\beta_0 - \beta_1) = 0 \quad (4.20)$$

Removing the brackets gives us:

$$\beta_0 + \beta_1 - \beta_0 + \beta_1 = 0 \quad (4.21)$$

This yields the equation:

$$2\beta_1 = 0 \quad (4.22)$$

Dividing both sides of the equation by 2, we get the null hypothesis that $\beta_1 = 0$.

Incidentally, if we had rescaled the contrast coding to be not ± 1 but $\pm 1/2$, the parameter β_1 would represent exactly the difference between the two means, and null hypothesis in equation (4.22) would have come out to be $\beta_1 = 0$. This is why it is sometimes better to recode the contrasts as $\pm 1/2$ rather than ± 1 . See [Schad et al. \(2020\)](#) for details; we will discuss this in the contrast coding chapter as well.

Let's load the data, set up the contrast coding, and fit the null versus the alternative models. We will fit varying intercept and varying slopes for subject and item, without correlations for items. We don't attempt to fit the so-called "maximal model" here with respect to the items random effects because we would get a singularity in the variance covariance matrix.

```
gg05e1 <- read.table("data/grodnergibsonE1crit.txt",
  header = TRUE)

gg05e1$so <- ifelse(gg05e1$condition == "objgap", 1,
  -1)
gg05e1$logrt <- log(gg05e1$rawRT)

library(lme4)
m0 <- lmer(logrt ~ 1 + (1 + so | subject) + (1 + so ||
  item), gg05e1)
m1 <- lmer(logrt ~ 1 + so + (1 + so | subject) + (1 +
  so || item), gg05e1)
```

Next, we compare the two models' log likelihoods. There is a function in the `lme4` package that achieves that: the `anova` function:

```
anova(m0, m1)
```

```
## refitting model(s) with ML (instead of REML)
```

```
## Data: gg05e1
## Models:
## m0: logrt ~ 1 + (1 + so | subject) + ((1 | item) + (0 + so | item))
## m1: logrt ~ 1 + so + (1 + so | subject) + ((1 | item) + (0 + so |
## m1:      item))
##      npar AIC BIC logLik deviance Chisq Df Pr(>Chisq)
## m0      7 707 739   -347      693
## m1      8 703 739   -343      687  6.15  1      0.013
```

You can confirm from the output that the `Chisq` value shown is minus two times the difference in log likelihood of the two models. The p-value is computed using the chi-squared distribution with one degree of freedom because in the two models the difference in the number of parameters is one:

```
round(pchisq(5.95, df = 1, lower.tail = FALSE), 3)
```

```
## [1] 0.015
```

It is common in the psycholinguistics literature to use the t-value from the linear mixed model output to conduct the hypothesis test on the slope:

```
summary(m1)$coefficients
```

```
##              Estimate Std. Error t value
## (Intercept)  5.88306    0.05176 113.669
## so           0.06202    0.02422   2.561
```

The more general method for hypothesis testing is the likelihood ratio test shown above.

One can also use the likelihood ratio test to evaluate whether a variance component should be included or not. For example, is the correlation parameter justified for the subjects random effects? Recall that we had a correlation of 0.58. Is this statistically significant? One can test this in the following way:

```

m1 <- lmer(logrt ~ 1 + so + (1 + so | subject) + (1 +
  so || item), gg05e1)
m1NoCorr <- lmer(logrt ~ 1 + so + (1 + so || subject) +
  (1 + so || item), gg05e1)
anova(m1, m1NoCorr)

## refitting model(s) with ML (instead of REML)

## Data: gg05e1
## Models:
## m1NoCorr: logrt ~ 1 + so + ((1 | subject) + (0 + so | subject)) + ((1 |
## m1NoCorr:      item) + (0 + so | item))
## m1: logrt ~ 1 + so + (1 + so | subject) + ((1 | item) + (0 + so |
## m1:      item))
##           npar AIC BIC logLik deviance Chisq Df
## m1NoCorr    7 710 741   -348      696
## m1           8 703 739   -343      687   8.7  1
##           Pr(>Chisq)
## m1NoCorr
## m1           0.0032

```

The test indicates that we can reject the null hypothesis that the correlation parameter is 0. We will return to this parameter in the chapter on simulation.

4.4 Exercises

4.4.1 Chinese relative clauses

Load the following two data-sets:

```

gibsonwu <- read.table("data/gibsonwucrit.txt", header = TRUE)
gibsonwu2 <- read.table("data/gibsonwu2012datarepeat.txt",
  header = TRUE)

```

The data are taken from two oexperiments that investigate (inter alia) the effect of relative clause type on reading time in Chinese. The data are from [Gibson and Wu \(2013\)](#) and [Vasishth et al. \(2013\)](#) respectively. The second data-set is a direct replication attempt of the first.

Chinese relative clauses are interesting theoretically because they are prenominal: the relative clause appears before the head noun. Figure 4.1 shows an example.

```
knitr::include_graphics("figures/chineserc.png")
```

The consequence of Chinese relative clauses being prenominal is that the distance between the gap in relative clause and the head noun is larger in subject relatives than object relatives (see Figure 4.1). [Hsiao and Gibson \(2003\)](#) have claimed that the larger distance in subject relatives leads to longer reading time at the head noun. Under this view, the prediction is that subject relatives are harder to process than object relatives. If this is true, this is interesting because in most other languages that have been studied, subject relatives are easier to process than object relatives; so Chinese will be a very unusual exception cross-linguistically.


The data provided are for the critical region (the head noun). The experiment method is self-paced reading, so we have reading times in milliseconds.

The research question is whether the difference in reading times between object and subject relative clauses is negative. For both data-sets, investigate this question by (a) fitting a paired t-test (by-subjects and by items), (b) fitting the most complex linear mixed model you can to the data and then interpreting the t-value, and (c) the likelihood ratio test. What can we conclude about the research question?

4.4.2 Agreement attraction in comprehension


Load the following data:

a. Subject relative

[GAP_i yaoqing fuhao de] guan

 GAP invite tycoon DE offici

‘The official who invited the tyc

b. Object relative

[fuhao yaoqing GAP_i de] guany

 tycoon invite GAP DE official

‘The official who the tycoon invi

FIGURE 4.1: An example of subjects vs. object relative clauses in Chinese.

```
datE1 <- read.table("data/dillonE1.txt", header = TRUE)
```

The data are taken from an experiment that investigate (inter alia) the effect of number similarity between a noun and the auxiliary verb in sentences like the following. There are two levels to a factor called Int(erference): low and high.

- low: The key to the cabinet *are* on the table
- high: The key to the *cabinets are* on the table

Here, in the condition marked *high*, the auxiliary verb *are* is predicted to be read faster than in the condition marked *low*, because the plural marking on the noun *cabinets* leads the reader to think that the sentence is grammatical. (Note that both sentences are ungrammatical.) This phenomenon, where the high condition is read faster than the low condition, is called **agreement attraction**.

The data provided are for the critical region (the auxiliary verb *are*). The experiment method is eyetracking; we have total reading times in milliseconds.

The research question is whether the difference in reading times between high and low conditions is negative.

- First, figure out which linear mixed model is appropriate for these data (varying intercepts only? varying intercepts and slopes? with or without correlations?).
- Then, carry out a statistical test using (a) the paired t-test (using the `t.test` function), (b) the t-test of the linear mixed model, and (c) the likelihood ratio test. What is your conclusion? Is there evidence for agreement attraction in the data?

4.4.3 The grammaticality illusion

Load the following data-sets:


```
english <- read.table("data/embeddingenglish.txt",  
  header = TRUE)  
dutch <- read.table("data/embeddingdutch.txt", header = TRUE)
```

In an offline accuracy rating study on English double center-embedding constructions, [Gibson and Thomas \(1999\)](#) found that grammatical constructions (e.g., example a below) were no less acceptable than ungrammatical constructions (e.g., example b) where a middle verb phrase (e.g., was cleaning every week) was missing.

- (a) The apartment that the maid who the service had sent over was cleaning every week was well decorated.
- (b) *The apartment that the maid who the service had sent over — was well decorated

Based on these results from English, [Gibson and Thomas \(1999\)](#) proposed that working-memory overload leads the comprehender to forget the prediction of the upcoming verb phrase (VP), which reduces working-memory load. This came to be known as the *VP-forgetting hypothesis*. The prediction is that in the word immediately following the final verb, the grammatical condition (which is coded as +1 in the data-frames) should be harder to read than the ungrammatical condition (which is coded as -1).

The data provided above test this hypothesis using self-paced reading for English ([Vasishth et al., 2011](#)), and for Dutch ([Frank et al., 2015](#)). The data provided are for the critical region (the noun phrase, labeled NP1, following the final verb). We have reading times in log milliseconds.

Is there support for the VP-forgetting hypothesis cross-linguistically, from English and Dutch?



5

Using simulation to understand your model

Data analysis is often taught as if the goal is to work out the p-value and make a decision: reject or fail to reject the null hypothesis. However, understanding the long-run properties of one's experiment design and statistical model under repeated sampling requires more work and thought. Specifically, it is important to understand (a) what one's model's power and Type I error properties are, and (b) whether the model we plan to fit to our data can, even in principle, recover the parameters in the model.

In order to study these properties of one's model, it is necessary to learn to simulate data that reflects our experimental design. Let's think about how to simulate data given a Latin-square 2 condition repeated measures design. We begin with our familiar running example, the [Grodner and Gibson \(2005\)](#) English relative clause data.

5.1 A reminder: The maximal linear mixed model

Recall the structure of the linear mixed model that can be used to fit the [Grodner and Gibson \(2005\)](#) data. We will discuss the so-called maximal model here—varying intercepts and slopes for subject and for item, with correlations—because that is the most general case.

In the model specification below, i indexes subjects, j items. The vector `so` has the sum contrast coding as usual: object relatives are coded as $+1/2$ and subject relatives as $-1/2$. We use this coding instead of ± 1 as before, because now the slope will reflect the

effect size rather than two times the effect size (see the hypothesis testing chapter).

Every row in the data-frame can be uniquely identified by the subject and item id, because this is a Latin square design and each subject sees exactly one instance of each item in a particular condition.

$$y_{ij} = \beta_0 + u_{0i} + w_{0j} + (\beta_1 + u_{1i} + w_{1j}) \times so_{ij} + \varepsilon_{ij} \quad (5.1)$$

where $\varepsilon_{ij} \sim \text{Normal}(0, \sigma)$ and

$$\Sigma_u = \begin{pmatrix} \sigma_{u0}^2 & \rho_u \sigma_{u0} \sigma_{u1} \\ \rho_u \sigma_{u0} \sigma_{u1} & \sigma_{u1}^2 \end{pmatrix} \quad \Sigma_w = \begin{pmatrix} \sigma_{w0}^2 & \rho_w \sigma_{w0} \sigma_{w1} \\ \rho_w \sigma_{w0} \sigma_{w1} & \sigma_{w1}^2 \end{pmatrix} \quad (5.2)$$

$$\begin{pmatrix} u_0 \\ u_1 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma_u \right), \quad \begin{pmatrix} w_0 \\ w_1 \end{pmatrix} \sim \mathcal{N} \left(\begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma_w \right) \quad (5.3)$$

β_0 and β_1 are the intercept and slope, representing the grand mean and the deviation from the grand mean in each condition. u are the subject level adjustments, and w the item level adjustments to the intercept and slope.

The above mathematical model expresses a generative process. In order to produce simulated data using the above process, we have to decide on some parameter values. We do this by estimating the parameters from the [Grodner and Gibson \(2005\)](#) study.

5.2 Obtain estimates from a previous study

First we load and prepare the relative clause data for data analysis.

```
gg05e1 <- read.table("data/grodnergibsonE1crit.txt",
  header = TRUE)
```

```
gg05e1$so <- ifelse(gg05e1$condition == "objgap", 1/2,
  -1/2)
gg05e1$logrt <- log(gg05e1$rawRT)
```

Next, fit the so-called maximal model. We will ignore the singularity warning as it won't affect us in our simulations.

```
library(lme4)
m<-lmer(logrt ~ so +
  (1+so|subject) +
  (1+so|item),
  data=gg05e1,
  ## "switch off" warnings:
  control=lmerControl(calc.derivs=FALSE))
```

The model summary shows that we can reject the null hypothesis of no difference in relative clauses:

```
summary(m)$coefficients
```

##	Estimate	Std. Error	t value
## (Intercept)	5.883	0.05202	113.082
## so	0.124	0.04932	2.515

Let's focus on that effect for our power analysis. What is the prospective power of detecting this effect *for a future study*? Note that we never compute power for an existing study—that is called post-hoc power and is a pointless quantity to compute because once the p-value is known, the power is just a transformation of the p-value (Hoenig and Heisey, 2001).

What we are doing below will *look* like post-hoc power because we are using existing data to compute power. However, what is crucially different in our approach is that (a) we remain unsure about the true effect, (b) we are making a statement about what

the power properties would be *if we ran the same study again*, with new subjects, but in the same environment (lab, etc.). We are not making any claim about the power properties of the *current* experiment; that ship has already sailed, the data are already at hand! Once the data are analyzed, it's too late to compute power for that particular data-set. A power analysis is only relevant for a design to be run in the future.

5.3 Decide on a range of plausible values of the effect size

Notice that the effect in milliseconds is relatively large, given the estimates from similar phenomena in reading studies in psycholinguistics (Jäger et al., 2017):

```
b0 <- summary(m)$coefficients[1, 1]
b1 <- summary(m)$coefficients[2, 1]
## effect estimate in log ms:
b1

## [1] 0.124

## effect estimate in ms:
exp(b0 + b1 * (0.5)) - exp(b0 + b1 * (-0.5))

## [1] 44.54
```

But the standard errors tell us that the effect could be as small or as large as the following values:

```
b1_stderr <- summary(m)$coefficients[2, 2]
lower <- b1 - 2 * b1_stderr
upper <- b1 + 2 * b1_stderr
lower
```

```
## [1] 0.02539
```

```
upper
```

```
## [1] 0.2227
```

The above range 0.03 and 0.22 arises because the range of plausible effect sizes is between $\hat{\beta}_1 \pm 2SE$ on the log ms scale.

On the ms scale, the range is:

```
exp(b0 + lower * (0.5)) - exp(b0 + lower * (-0.5))
```

```
## [1] 9.113
```

```
exp(b0 + upper * (0.5)) - exp(b0 + upper * (-0.5))
```

```
## [1] 80.09
```

On the ms scale we see that that's a *lot* of uncertainty in the effect size! With some experience, you will come to recognize that such a wide confidence bound is a sign of low power. We will just establish the prospective power properties of this study in a minute.

We can take the above uncertainty of the $\hat{\beta}_1$ estimator into account (on the log ms scale—remember that the model is based on log rt) by assuming that the effect has the following uncertainty on the log ms scale:

$$\beta_1 \sim \text{Normal}(0.12, 0.05) \quad (5.4)$$

Here, we are doing something that is, strictly speaking, Bayesian in thinking. We are describing our uncertainty about the true effect from the best estimate we have—existing data. To talk about the uncertainty, we are (ab)using the 95% confidence interval (treating it like its telling us the range of plausible values). Recall that strictly speaking, in the frequentist paradigm, one cannot talk

about the probability distribution of the effect size—in frequentist theory, the true value of the parameter is a point value, it has no distribution. The range $\hat{\beta}_1 \pm 2 \times SE$ refers to the estimated mean of the sampling distribution of the sample means, and to the standard deviation of this sampling distribution. Thus, strictly speaking, this range does not reflect our uncertainty about the true parameter's value. Having said this, we are going to use the effect estimates from our model fit as a starting point for our power analysis because this is the best information we have so far about the English relative clause design.

5.4 Extract parameter estimates

Next, in preparation for the power analysis, we extract all the parameter estimates from the model we have fit above. The parameters are:

- The two fixed effects (the β parameters)
- The residuals' standard deviation
- The standard deviations of the subject intercept and slope adjustments, and the corresponding correlation matrix.
- The standard deviations of the item intercept and slope adjustments, and the corresponding correlation matrix.

The correlation matrices and the subject/item random effects standard deviations are used to assemble the variance covariance matrix; this is done using the `sdcor2cov` function from the `SIN` package; recall the discussion in chapter 1. For the variance covariance matrix for items random effects, we use an intermediate value of 0.5 for the correlation parameter because the linear mixed model was unable to estimate the parameter.

```
## extract parameter estimates:
beta <- round(summary(m)$coefficients[, 1], 4)
sigma_e <- round(attr(VarCorr(m), "sc"), 2)
```



```

subj_ranefsd <- round(attr(VarCorr(m)$subject, "stddev"),
  4)
subj_ranefcorr <- round(attr(VarCorr(m)$subject, "corr"),
  1)

## assemble variance-covariance matrix for subjects:
Sigma_u <- SIN::sdcor2cov(stddev = subj_ranefsd, corr = subj_ranefcorr)
## check that the matrix can be inverted:
solve(Sigma_u)

```

```

##           (Intercept)      so
## (Intercept)      15.46 -13.31
## so              -13.31  31.82

```

```

item_ranefsd <- round(attr(VarCorr(m)$item, "stddev"),
  4)
item_ranefcorr <- round(attr(VarCorr(m)$item, "corr"),
  1)

## assemble variance matrix for items: ## this won't
## work:
## Sigma_w<-SIN::sdcor2cov(stddev=item_ranefsd,
## corr=item_ranefcorr) solve(Sigma_w) choose some
## intermediate values for correlations:
corr_matrix <- (diag(2) + matrix(rep(1, 4), ncol = 2))/2

Sigma_w <- SIN::sdcor2cov(stddev = item_ranefsd, corr = corr_matrix)
## matrix inverts:
solve(Sigma_w)

```

```

##           [,1]    [,2]
## [1,]   774.2 -181.3
## [2,] -181.3   169.9

```

5.5 Define a function for generating data

Next, we define a function that generates repeated measures data given the parameter estimates. The basic idea here is the following.

- First, create a data-frame that represents a Latin-square design.
- Then, given the condition id, and the subject and item ids in each row of the data frame, generate data row-by-row.

We explain these steps next.

5.5.1 Generate a Latin-square design

First, consider how one can create a Latin-square design. Suppose we have four items and four subjects. For such an experiment, we would create two groups, g1 and g2, with the following layout.

```
nitem <- 4
nsubj <- 4

g1 <- data.frame(item = 1:nitem, cond = rep(c("a",
      "b"), nitem/2))
g2 <- data.frame(item = 1:nitem, cond = rep(c("b",
      "a"), nitem/2))
g1
```

```
##   item cond
## 1     1    a
## 2     2    b
## 3     3    a
## 4     4    b
```

```
g2
```

```
##   item cond
## 1     1    b
```

```
## 2    2    a
## 3    3    b
## 4    4    a
```

Half the total number of subjects will be assigned to group 1 and half to group 2:

```
## assemble data frame:
gp1 <- g1[rep(seq_len(nrow(g1)), nsubj/2), ]
gp2 <- g2[rep(seq_len(nrow(g2)), nsubj/2), ]

simdat <- rbind(gp1, gp2)

## add subject column:
simdat$subj <- rep(1:nsubj, each = nitem)
```

Finally, the contrast coding for each row in the data-frame is set up:

```
## add contrast coding:
simdat$so <- ifelse(simdat$cond == "a", -1/2, 1/2)
```

5.5.2 Generate data row-by-row

Then, we proceed row-by-row in this data frame, and generate data for each subject, item, and condition. For example, the first row of our simulated data-set has subject id:

```
simdat[1, ]$subj
```

```
## [1] 1
```

Similarly, the first row has item id:

```
simdat[1, ]$item
```

```
## [1] 1
```

The first row's condition coding is:

```
simdat[1, ]$so
```

```
## [1] -0.5
```

These three pieces of information are what we need to generate data for the first row. Recall that the model for subject i and item j in condition `so` is

$$\beta_0 + u_{0i} + w_{0j} + (\beta + 1 + u_{1i} + w_{1j}) \times so + \varepsilon \text{ where } \varepsilon \sim N(0, \sigma) \quad (5.5)$$

The terms `u0`, `w0`, and `u1`, `w1`, which are the adjustments to the intercepts and slopes by subject and item, are stored in two matrices that are generated randomly each time that we simulate new subjects/items. Recall from chapter 1 how bivariate data are generated. The intercept and slope adjustments will be generated using the `mvrnorm` function in the `MASS` library. For example, given the variance covariance matrix for subjects `Sigma_u` that we created above, the subject random effects (intercept and slope adjustments) for 10 subjects can be generated in a matrix as follows:

```
library(MASS)
u <- mvrnorm(n = 10, mu = c(0, 0), Sigma = Sigma_u)
u
```

```
##      (Intercept)      so
## [1,]      0.1059 -0.06177
## [2,]      0.4790 -0.13362
## [3,]      0.3538  0.26910
## [4,]      0.1804  0.22173
## [5,]     -0.1610  0.34222
## [6,]      0.2322  0.08877
## [7,]      0.3389  0.16131
## [8,]     -0.2254  0.12916
```

```
## [9,]      -0.1795 -0.20251
## [10,]      0.6060  0.40070
```

Each row in this matrix is the intercept and slope adjustment for a subject; the row number indexes the subject id. For example, subject 1's intercept adjustment is:

```
u[1, 1]
```

```
## (Intercept)
##      0.1059
```

Subject 1's slope adjustment is:

```
u[1, 2]
```

```
##      so
## -0.06177
```

Analogously to the subject random effects matrix, a matrix for items random effects is also generated. As an example, we generate simulated random effects for 10 items:

```
w <- mvrnorm(n = 10, mu = c(0, 0), Sigma = Sigma_w)
w
```

```
##           [,1]      [,2]
## [1,] -0.028129 -0.0476951
## [2,]  0.019630  0.1601535
## [3,]  0.060869 -0.0519072
## [4,] -0.057216 -0.0980981
## [5,]  0.025503  0.1211837
## [6,]  0.047698 -0.0810234
## [7,]  0.014671  0.0206600
## [8,] -0.005208 -0.0001712
## [9,]  0.097832  0.0729380
## [10,] -0.013452 -0.0386293
```

Now, to generate simulated data for subject 1 and item 1 for object relatives, we simply need to run this line of code:

```
rlnorm(1, beta[1] + u[1, 1] + w[1, 1] + (beta[2] +
  u[1, 2] + w[1, 2]) * (0.5), sigma_e)
```

```
## [1] 398.3
```

For subject 2, all that would change is the subject id: instead of `u[1,1]`, and `u[1,2]` we would write `u[2,1]` and `u[2,2]`. The for-loop below works through the `simdat` data-frame row by row, looks up the subject id, the item id for that row, and the condition coding for that row, and fills in the simulated reading time using the above code. This is how the simulated data are generated.

Here is the complete function for generating simulated data. It uses all the bits of code we discussed above.

```
library(MASS)
## assumes that no. of subjects and
## no. of items is divisible by 2.
gen_sim_lnorm2<-function(nitem=16,
  nsubj=42,
  beta=NULL,
  Sigma_u=NULL, # subject vcov matrix
  Sigma_w=NULL, # item vcov matrix
  sigma_e=NULL){
  ## prepare data frame for a two-condition latin square:
  g1<-data.frame(item=1:nitem,
    cond=rep(c("a","b"),nitem/2))
  g2<-data.frame(item=1:nitem,
    cond=rep(c("b","a"),nitem/2))

  ## assemble data frame:
  gp1<-g1[rep(seq_len(nrow(g1)),
    nsubj/2),]
```

```

gp2<-g2[rep(seq_len(nrow(g2)),
             nsubj/2),]

simdat<-rbind(gp1,gp2)

## add subject column:
simdat$subj<-rep(1:nsubj,each=nitem)

## add contrast coding:
simdat$so<-ifelse(simdat$cond=="a",-1/2,1/2)

## subject random effects:
u<-mvrnorm(n=length(unique(simdat$subj)),
           mu=c(0,0),Sigma=Sigma_u)

## item random effects
w<-mvrnorm(n=length(unique(simdat$item)),
           mu=c(0,0),Sigma=Sigma_w)

## generate data row by row:
N<-dim(simdat)[1]
rt<-rep(NA,N)
for(i in 1:N){
  rt[i] <- rlnorm(1,beta[1] +
                u[simdat[i,]$subj,1] +
                w[simdat[i,]$item,1] +
                (beta[2]+u[simdat[i,]$subj,2]+
                 w[simdat[i,]$item,2])*simdat$so[i],
                sigma_e)
}
simdat$rt<-rt
simdat$subj<-factor(simdat$subj)
simdat$item<-factor(simdat$item)
simdat}

```

Let's generate some simulated data and check what the data look like:

```
dat <- gen_sim_lnorm2(nitem = 16, nsubj = 42, beta = beta,
  Sigma_u = Sigma_u, Sigma_w = Sigma_w, sigma_e = sigma_e)
```

The data have the expected structure:

```
## fully crossed subjects and items:
head(t(xtabs(~subj + item, dat)))
```

```
##      subj
## item 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19
##      1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      3 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      4 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      5 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      6 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      subj
## item 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35
##      1 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      2 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      3 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      4 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      5 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      6 1 1 1 1 1 1 1 1 1 1 1 1 1 1 1
##      subj
## item 36 37 38 39 40 41 42
##      1 1 1 1 1 1 1
##      2 1 1 1 1 1 1
##      3 1 1 1 1 1 1
##      4 1 1 1 1 1 1
##      5 1 1 1 1 1 1
##      6 1 1 1 1 1 1
```



```
## 8 measurements per condition:
head(t(xtabs(~subj + cond, dat)))
```

```
##      subj
## cond 1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19
##      a 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
##      b 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
##      subj
## cond 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35
##      a 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
##      b 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8 8
##      subj
## cond 36 37 38 39 40 41 42
##      a 8 8 8 8 8 8 8
##      b 8 8 8 8 8 8 8
```

```
## contrast coding check:
xtabs(~so + cond, dat)
```

```
##      cond
## so      a    b
##  -0.5 336    0
##   0.5   0 336
```

```
## condition b is slower than a:
round(with(dat, tapply(rt, cond, mean)))
```

```
##      a    b
## 356 431
```

Everything checks out.

5.6 Repeated generation of data to compute power

With the function for simulating data ready, we are now able to repeatedly generate simulated data. Next, we generate data 100 times, fit a linear mixed model each time, and extract the t-value from each linear mixed model fit. The vector t-values is then used to compute power: we simply look at the proportion of absolute t-values that exceed the value 2.

```
nsim <- 100
sotval <- rep(NA, nsim)

for (i in 1:nsim) {
  # generate sim data:
  dat <- gen_sim_lnorm2(nitem = 16, nsubj = 42, beta = beta,
    Sigma_u = Sigma_u, Sigma_w = Sigma_w, sigma_e = sigma_e)

  ## fit model to sim data:
  m <- lmer(log(rt) ~ so + (1 + so | subj) + (1 +
    so | item), dat, control = lmerControl(calc.derivs = FALSE))
  ## extract the t-value
  sotval[i] <- summary(m)$coefficients[2, 3]
}

mean(abs(sotval) > 2)
```

```
[1] 0.65
```

The following computation will take a lot of time because we are generating and fitting data 3×100 times.

Here, we will assume that the true effect has the following range of plausible values:

```
## lower bound:
0.12 - 2 * 0.05
```

```
## [1] 0.02
```

```
## mean
0.12
```

```
## [1] 0.12
```

```
## upper bound
0.12 + 2 * 0.05
```

```
## [1] 0.22
```

We will run two for-loops now: the first for-loop selects one of the plausible values as the value for the slope, and then the second for-loop runs 100 simulations to compute power for that effect size.

This time, instead of ignoring convergence problems, we can record the proportion of times that we get a convergence failure or problem, and we can discard that result:

```
nsim <- 100
## effect size possibilities:
b1_est <- c(0.02, 0.12, 0.22)
sotvals <- matrix(rep(NA, nsim * length(b1_est)), ncol = nsim)
failed <- matrix(rep(0, nsim * length(b1_est)), ncol = nsim)
for (j in 1:length(b1_est)) {
  for (i in 1:nsim) {
    beta[2] <- b1_est[j]
    dat_sim <- gen_sim_lnorm2(nitem = 16, nsubj = 42,
      beta = beta, Sigma_u = Sigma_u, Sigma_w = Sigma_w,
      sigma_e = sigma_e)

    ## no correlations estimated to minimize convergence
    ## problems: analysis done after log-transforming:
    m <- lmer(log(rt) ~ so + (so || subj) + (so ||
      item), data = dat_sim)
    ## ignore failed trials
```

```

    if (any(grepl("failed to converge", m@optinfo$conv$lme4$messages))) {
      failed[j, i] <- 1
    } else {
      sotvals[j, i] <- summary(m)$coefficients[2,
        3]
    }
  }
}
## proportion of convergence failures:
rowMeans(failed)

```

Power can now be computed for each effect size:

```

pow <- rep(NA, length(b1_est))
for (k in 1:length(b1_est)) {
  pow[k] <- mean(abs(sotvals[k, ]) > 2, na.rm = TRUE)
}

pow

```

```
## [1] 0.06122 0.72000 0.98990
```

Notice that there is a lot of uncertainty about the power estimate here!

Recall that power is a **function** of

- effect size
- standard deviation(s); in linear mixed models, these are all the variance components and the correlations
- sample size (numbers of subjects and items)

In papers, you will often see text like “power was x%”. This statement reflects a misunderstanding; power is best plotted as a function of (a subset of) these variables.

In the discussion above, we display a range of power estimates; this range reflects our uncertainty about the power estimate as a

function of the plausible effect sizes. Often (as here) this uncertainty will be very high! I.e., given what one knows so far, it may be difficult to pin down what the assumed effect size etc., should be, and that makes it hard to give a precise range for your power estimate.

5.7 What you can now do

Given the above code and workflow, you can now figure out how many subjects you might need to achieve 80% power, assuming a certain effect size (or a range of effect sizes as above) and assuming some specific values for the standard deviations and variance covariance matrices.

You can also study Type I error properties of your model as a function of whether the model is a maximal model or a varying intercepts only model.

Example: Compute power as a function of effect size and sample size. Note that the number of subjects has to be even, otherwise the simulation code will fail! One could put in a test for this in the code: if the number of subjects is divisible by 2, the modulo function should return 0:

```
10%%2
```

```
## [1] 0
```

```
11%%2
```

```
## [1] 1
```

```
## define a function for computing power (as a  
## function of effect size and subject sample size:
```

```

compute_power <- function(b = NULL, nsubjects = 28) {
  if (nsubjects%%2 != 0) {
    stop("No. of subjects must be divisible by 2.")
  }
  nsim <- 100
  sotvals <- rep(NA, nsim)
  failed <- rep(0, nsim)
  for (i in 1:nsim) {
    beta[2] <- b
    dat_sim <- gen_sim_lnorm2(nitem = 24, nsubj = nsubjects,
      beta = beta, Sigma_u = Sigma_u, Sigma_w = Sigma_w,
      sigma_e = sigma_e)

    ## no correlations estimated to avoid convergence
    ## problems: analysis done after log-transforming:
    m <- lmer(log(rt) ~ so + (so || subj) + (so ||
      item), data = dat_sim)
    ## ignore failed trials
    if (any(grepl("failed to converge", m@optinfo$conv$lme4$messages))) {
      failed[i] <- 1
    } else {
      sotvals[i] <- summary(m)$coefficients[2,
        3]
    }
  }

  ## proportion of convergence failures and power:
  paste(print("Prop. of convergence failures:"),
    mean(failed), sep = " ")
  paste(print("Power:"), mean(abs(sotvals) > 2, na.rm = TRUE),
    sep = " ")
}

```

```

## usage: this will halt function with an error
## message compute_power(b=0.03,nsubjects=29)

```

```
compute_power(b = 0.03, nsubjects = 28)
```

```
[1] "Prop. of convergence failures:" [1] "Power:" [1] "Power:  
0.0707070707070707"
```

5.8 Exercises

5.8.1 Drawing a power curve given a range of effect sizes

Use the simulation code as provided to compute a power function for effects sizes for the relative clause effect ranging from 0.025, 0.05, 0.10, and 0.15, given that you have 16 items and 42 participants.

5.8.2 Power and log-transformation

Modify the simulation code to generate not log-normally distributed data, but normally distributed data. Refit the [Grodner and Gibson \(2005\)](#) data using raw reading times (i.e., do not log-transform them), and then use the parameter estimates from the data to compute a power function for effects sizes for the relative clause effect ranging from 10, 30, 60, 80 ms, given that you have 16 items and 42 participants. Compare your power curve with that of Part 1.

5.8.3 Evaluating models by generating simulated data

Generate data from the simulation function assuming a log-normal likelihood and then generate data from the function you wrote in Part 2 that assumes a normal likelihood. Compare the distributions of the two sets of simulated data to the observed distributions. Which simulation code produces more realistic data, and why?

5.8.4 Using simulation to check parameter recovery

Check whether the simulation code you wrote assuming a normal likelihood can recover the parameters.

5.8.5 Sample size calculations using simulation

Load the data-set shown below:

```
gibsonwu <- read.table("data/gibsonwucrit.txt")
```

Use simulation to determine how many subjects you would need to achieve power of 80%, given 16 items, and an effect size of 0.02 on the log ms scale. Draw a power curve: on the x-axis show the number of subjects, and on the y-axis the estimated power. Now draw two further curves, one for an effect size of 0.05 and another for an effect size of 0.10. This gives you a power curve, taking the uncertainty in the effect size into account.

6

Important distributions

These distributions are used quite frequently in Bayesian data analyses, especially in psychology and linguistics applications. The Binomial and Poisson are discrete distributions, the rest are continuous. Each distribution comes with a family of **d-p-q-r** functions in R which allow us to compute the PDF/PMF, the CDF, the inverse CDF, and to generate random data. For example, the normal distribution's PDF is **dnorm**; the CDF and the inverse CDF are **pnorm** and **qnorm** respectively; and random data can be generated using **rnorm**.

The table below is adapted from https://github.com/wzchen/probability_cheatsheet, which is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License.

to-do: check that the notation is consistent with the main text's.

Distribution	PMF/PDF and Support	Expected Value	Variance
Binomial <i>Binomial</i> (n, θ)	$P(X = k) = \binom{n}{k} \theta^k (1 - \theta)^{n-k}$ $k \in \{0, 1, 2, \dots, n\}$	$n\theta$	$n\theta(1 - \theta)$
Poisson <i>Pois</i> (λ)	$P(X = k) = \frac{e^{-\lambda} \lambda^k}{k!}$ $k \in \{0, 1, 2, \dots\}$	λ	λ
Uniform <i>Unif</i> (a, b)	$f(x) = \frac{1}{b-a}$ $x \in (a, b)$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$
Normal <i>Normal</i> (μ, σ)	$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{(2\sigma^2)}}$ $x \in (-\infty, \infty)$	$\mu = \frac{\sum_{i=1}^n x_i}{n}$	$\sigma^2 = \frac{\sum_{i=1}^n (x_i - \bar{x})^2}{n}$
Log-Normal <i>LogNormal</i> (μ, σ)	$\frac{1}{x\sigma\sqrt{2\pi}} e^{-(\log x - \mu)^2 / (2\sigma^2)}$ $x \in (0, \infty)$	$\theta = e^{\mu + \sigma^2/2}$	$\theta^2(e^{\sigma^2} - 1)$
Beta <i>Beta</i> (a, b)	$f(x) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} x^{a-1} (1-x)^{b-1}$ $x \in (0, 1)$	$\mu = \frac{a}{a+b}$	$\frac{\mu(1-\mu)}{(a+b+1)}$
Exponential <i>Exp</i> (λ)	$f(x) = \lambda e^{-\lambda x}$ $x \in (0, \infty)$	$\frac{1}{\lambda}$	$\frac{1}{\lambda^2}$
Gamma <i>Gamma</i> (a, λ)	$f(x) = \frac{1}{\Gamma(a)} (\lambda x)^a e^{-\lambda x} \frac{1}{x}$ $x \in (0, \infty)$	$\frac{a}{\lambda}$	$\frac{a}{\lambda^2}$
Student- t $t(n)$ Cauchy is $t(1)$	$\frac{\Gamma((n+1)/2)}{\sqrt{n\pi}\Gamma(n/2)} (1 + x^2/n)^{-(n+1)/2}$ $x \in (-\infty, \infty)$	0 if $n > 1$	$\frac{n}{n-2}$ if $n > 2$

Bibliography

- Barr, D. J., Levy, R., Scheepers, C., and Tily, H. J. (2013). Random effects structure for confirmatory hypothesis testing: Keep it maximal. *Journal of Memory and Language*, 68(3):255–278.
- Bates, D. M., Kliegl, R., Vasishth, S., and Baayen, H. (2015). Parsimonious mixed models. Unpublished manuscript.
- Benjamin, D. J., Berger, J. O., Johannesson, M., Nosek, B. A., Wagenmakers, E.-J., Berk, R., Bollen, K. A., Brembs, B., Brown, L., Camerer, C., et al. (2018). Redefine statistical significance. *Nature Human Behaviour*, 2(1):6.
- Blitzstein, J. K. and Hwang, J. (2014). *Introduction to probability*. Chapman and Hall/CRC.
- Button, K. S., Ioannidis, J. P., Mokrysz, C., Nosek, B. A., Flint, J., Robinson, E. S., and Munafò, M. R. (2013). Power failure: why small sample size undermines the reliability of neuroscience. *Nature Reviews Neuroscience*, 14(5):365–376.
- Cumming, G. (2014). The new statistics: Why and how. *Psychological science*, 25(1):7–29.
- Fox, J. (2009). *A mathematical primer for social statistics*. Number 159. Sage.
- Frank, S. L., Trompenaars, T., and Vasishth, S. (2015). Cross-linguistic differences in processing double-embedded relative clauses: Working-memory constraints or language statistics? *Cognitive Science*, 40:554–578.
- Gelman, A., Carlin, J. B., Stern, H. S., Dunson, D. B., Vehtari, A., and Rubin, D. B. (2014). *Bayesian Data Analysis*. Chapman and Hall/CRC, Boca Raton, FL, third edition.
- Gibson, E. and Thomas, J. (1999). Memory limitations and structural

- forgetting: The perception of complex ungrammatical sentences as grammatical. *Language and Cognitive Processes*, 14(3):225–248.
- Gibson, E. and Wu, H.-H. I. (2013). Processing Chinese relative clauses in context. *Language and Cognitive Processes*, 28(1-2):125–155.
- Grodner, D. and Gibson, E. (2005). Consequences of the serial nature of linguistic input. *Cognitive Science*, 29:261–290.
- Hedges, L. V. (1984). Estimation of effect size under nonrandom sampling: The effects of censoring studies yielding statistically insignificant mean differences. *Journal of Educational Statistics*, 9(1):61–85.
- Hoening, J. M. and Heisey, D. M. (2001). The abuse of power: The pervasive fallacy of power calculations for data analysis. *The American Statistician*, 55:19–24.
- Hsiao, F. P.-F. and Gibson, E. (2003). Processing relative clauses in Chinese. *Cognition*, 90:3–27.
- Ioannidis, J. P. (2008). Why most discovered true associations are inflated. *Epidemiology*, 19(5):640–648.
- Jäger, L. A., Engelmann, F., and Vasishth, S. (2017). Similarity-based interference in sentence comprehension: Literature review and Bayesian meta-analysis. *Journal of Memory and Language*, 94:316–339.
- Johnson, K. (2011). *Quantitative methods in linguistics*. John Wiley & Sons.
- Kerns, G. J. (2010). *Introduction to Probability and Statistics Using R*.
- Lane, D. M. and Dunlap, W. P. (1978). Estimating effect size: Bias resulting from the significance criterion in editorial decisions. *British Journal of Mathematical and Statistical Psychology*, 31(2):107–112.
- Matuschek, H., Kliegl, R., Vasishth, S., Baayen, R. H., and Bates, D. M. (2017). Balancing Type I Error and Power in Linear Mixed Models. *Journal of Memory and Language*, 94:305–315.
- Miller, I. and Miller, M. (2004). *John E. Freund’s Mathematical Statistics with Applications*. Prentice Hall.
- Morin, D. J. (2016). *Probability: For the Enthusiastic Beginner*. Createspace Independent Publishing Platform.

- Pocock, S. J. (2013). *Clinical trials: A practical approach*. John Wiley & Sons.
- R Core Team (2019). *R: A Language and Environment for Statistical Computing*. R Foundation for Statistical Computing, Vienna, Austria.
- Rice, J. A. (1995). *Mathematical statistics and data analysis*. Duxbury press Belmont, CA.
- Schad, D. J., Vasishth, S., Hohenstein, S., and Kliegl, R. (2020). How to capitalize on a priori contrasts in linear (mixed) models: A tutorial. *Journal of Memory and Language*, 110.
- Vasishth, S., Chen, Z., Li, Q., and Guo, G. (2013). Processing Chinese relative clauses: Evidence for the subject-relative advantage. *PLoS ONE*, 8(10):1–14.
- Vasishth, S., Mertzen, D., Jäger, L. A., and Gelman, A. (2018). The statistical significance filter leads to overoptimistic expectations of replicability. *Journal of Memory and Language*, 103:151–175.
- Vasishth, S., Suckow, K., Lewis, R. L., and Kern, S. (2011). Short-term forgetting in sentence comprehension: Crosslinguistic evidence from head-final structures. *Language and Cognitive Processes*, 25:533–567.

