

FLORIAN HARTIG

ESSENTIAL STATISTICS

Contents

<i>1</i>	<i>Introduction</i>	<i>3</i>
<i>2</i>	<i>Data, descriptive statistics and visualization</i>	<i>5</i>
<i>3</i>	<i>Inferential statistics</i>	<i>11</i>
<i>4</i>	<i>Predictive statistics - machine learning</i>	<i>23</i>
<i>5</i>	<i>Design of experiments</i>	<i>25</i>
<i>6</i>	<i>Good to knows and further reading</i>	<i>28</i>
<i>7</i>	<i>Bibliography</i>	<i>30</i>

1

Introduction

1.1 Purpose and intended audience

This document provides a short introduction to the field of statistics, as well as to statistical analyses most likely to be encountered in elementary experimental and observational situations. It is intended as a primer and will not replace a more thorough lecture or textbook. Regarding the latter, I provide recommendations in the final section [6.2](#).

1.2 Topics of statistics and data science

Statistics, or more broadly data science, deals with the visualization, summary and interpretation of data. This primer provides an introduction to the four most important pillars of statistical methods for a quantitative researcher:

Descriptive statistics include summary statistics such as mean and median as well as the various options to visualize data.

Inferential statistics deals with testing hypothesis and fitting models, typically based on a statistical model (the data-generating process).

Predictive statistics and machine learning deals with deriving predictions from data, in particular from "big data".

Experimental design covers all aspects of generating data, in particular questions such as "which variables need to be measured?", "how many data do we need?", "how should we optimally vary the variables of interest in an experiment"?

1.3 The R environment for statistical computing

For the better or worse, the times when statistics was done with pen, paper and a calculator are over. Statistical analysis nowadays happens on the computer, and a number of software environments exist to do so. In the ecological sciences, R is the de-facto standard

The difference between inferential and predictive statistics is that inferential statistics is typically concerned with a decision about a model or assumption (is the hypothesis true ...), while predictive methods deal with making predictions from data, often without the intermediate step of modelling a "data-generating process".

for statistical analysis. R is open-source, free, and has a very larger user base, specially in the environmental sciences, that contribute packages for specialised ecological and environmental analysis.

In this primer, all examples will be calculated with R; however, I will give no further introduction to R. If you need one, follow [this link](#) for an introduction, including help on how to install the software. If you start using R for the first time, I highly recommend to use R together with RStudio. Follow the link above to get further information

R is a script-based language, which means that you communicate with the computer not by clicking on buttons, but by writing commands either directly in the R console, or first in a text file and then sending it to the R console, which is then evaluating your commands. If you aren't used to this kind of approach yet, it may take a short while to get used to this, but once you are used to it, you will notice how advantageous it is to have all steps of your analysis listed in a text file, being able to repeat everything at any moment.

2

Data, descriptive statistics and visualization

Descriptive statistics deals with the summary and visualization of data. Before we speak about those methods, it will be useful to shortly discuss what data is and what properties it may have.

We will speak about how to create data in the later chapter 5 on experimental design.

2.1 Data

You can think of this situation as a spreadsheet where the columns are the variables and the rows are the observations.

Of course, there are other data structures, but this is the most common one.

A typical dataset consists of multiple observations of number of variables (e.g. temperature, precipitation, growth).

The response variable is the variable for which we try to understand how it responds to other factors.

Usually, this data will contain one variable on which we want to focus. We call this variable the response variable (sometimes also the dependent variable or outcome variable), because we are interested if and how this variable of interest varies (responds, depends) when something else changes. The variables that affect the response could be an environmental factor (e.g. temperature), it could be an experimental treatment (fertilized vs. non fertilized), or anything else. Those other variables that affect our response are called predictor variables (synonymous terms are explanatory variables, covariates or independent variables).

The predictor variables are those that affect the response.

The most common case is that the response variable is a single value (e.g. a number or a categorical outcome), and we will concentrate on this case. However, there are cases when the response has more than one dimension, or when we are interested in the change of several response variables at a time. If the response has several dimensions, we have to use multivariate analysis methods, which are not explained here. The analysis of such data is known as multivariate statistics. We will not cover this here, but if you need it some further links are [here](#).

Another important distinction is the type of variables. Independent of whether we are speaking about the response or the predictor, we distinguish:

Variables can be continuous, discrete or categorical. Categorical variables can be ordered, unordered, or binary.

- Continuous numeric variables (ordered and dense), e.g. temperature
- Discrete numeric variables (ordered, but discrete), e.g. count data

- Categorical variables (e.g. a fixed set of options such as red, green blue), which can further be divided into
 - Unordered categorical variables (Nominal) such as red, green, blue
 - Binary (Dichotomous) variables (dead / survived)
 - Ordered categorical variables (small, medium, large)

Experience shows that there is certain tendency of beginners to use categorical variables for things that are actually continuous, e.g. by coding body weight of animals into light, medium, heavy. The justification stated is often that this avoids the measurement uncertainty. In short: it doesn't, it just creates more problems. Don't use categorical variables for things that can also be recorded numerically!

Don't use categorical variables for things that can also be recorded numerically!

In R:

To represent data, R has a basic data structure, the `data.frame`. A data frame is like a spread sheet, with columns, and each column has can have a different type. Possibilities are

- integer - what it says
- numeric - continuous number (float)
- boolean - true / false
- factor - normally unordered, i.e. red, green, blue. Can also be made ordered (small, medium, large), although it is then often better to code this as an integer

Also, although not really a data type, a common case is the absence of an observation for a particular variable. This is typically code by "NA". Similar, but not the same is "NaN" (not a number), which occurs as the result of a calculation that cannot be performed.

To read in data as a `data.frame` in R-studio, go to the upper right part, "Import Dataset". For more explanation, see <http://biometry.github.io/APES/R/R20-DataStructures.html>

It is very important that you check that your columns have the right type after reading in the data. To do so, type in `str(TheNameOfMyData)`, and check the type of the columns

2.2 Summary statistics

Summary statistics are numerical calculations that summarize a dataset. They are used to display the properties of a dataset in a more compact way.

A simple summary statistics you should be familiar with is the mean

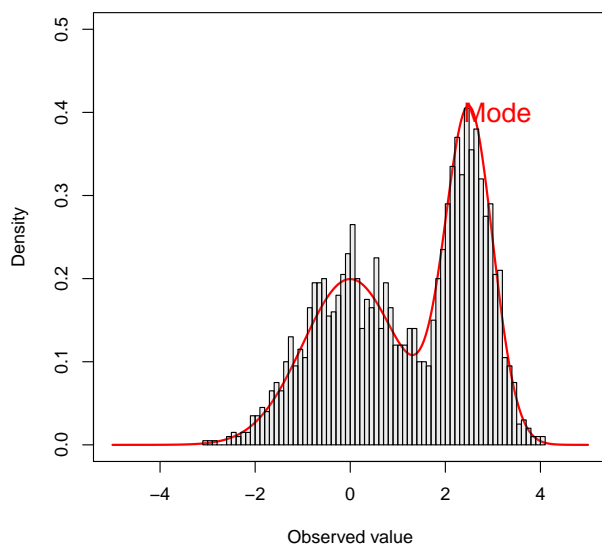
Summarizing a single continuous variable

Imagine we have a single, continuous variable with a large number of observations - hence, we have a distribution of values (see below).

How can we summarize the properties of this distribution. Some basic properties would be the minimum and the maximum, as well as the mode (the maximum of the distribution, i.e. value with the highest density of observations). There are two further central summaries that are always calculated - moments and quantiles.

The n -th moment is defined as the average across the distribution, taking each observed value to the power of n .

Figure 2.1: bla



$$\langle x^2 \rangle_d \quad (2.1)$$

where $\langle \rangle$ denotes the average, and d the distribution. The first four moments are called the mean (average value), the standard deviation (measure of spread), the skewness (measure of asymmetry in the distribution) and the kurtosis.

The second central quantity for describing distributions are the quantiles. If we have a distribution such as the figure above, we can ask ourself: which is the value of the variable that divides the data in half, so that half of the observed data are lower, and half are higher than this value? This point is called the median, and also the 0.5 quantile. More general, the 0.x quantile is the value at which a fraction of 0.x of the data is smaller.

Half of the data are lower, and half of the data are higher than the median, the 0.5 quantile

Correlatins - summarizing the dependence between continous variables

Correlation measures the dependence between continous variables.

Add some example of pearson correlation here

Mention that there are other correlations, i.e. rank correlation

Contingency tables - summarizing disrecte outcomes of several variables

Contingency tables measure

add an example of contingency tables here

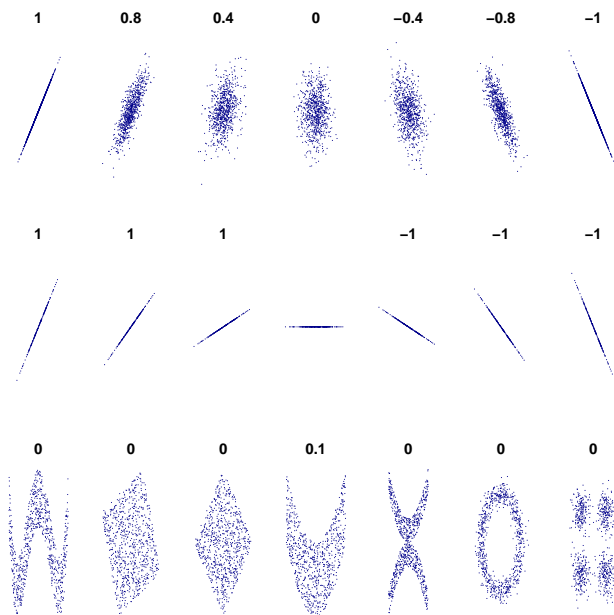


Figure 2.2: bla

In R:

For calculating descriptive statistics in R, see [here](#)

2.3 Visualization

Summary statistics are useful, but also dangerous. A famous example is Anscombe's Quartet, a hypothetical dataset of four observations that are identical in classical summary statistics such as mean, variance, correlation, regression line, etc. (?).

Type `?anscombe` in R to see the code to create these plots and to calculate the statistical properties of the datasets

Principles of visualization

Clean graphs

Truthfull

Not more than necessary

Basic graphs

Line plots continuous continuous, ordered

Scatter plots continuous continuous, ordered

Bar plots

Box plots

In R:

There are number of excellent introductions to graphics with R, so that I won't bother to provide

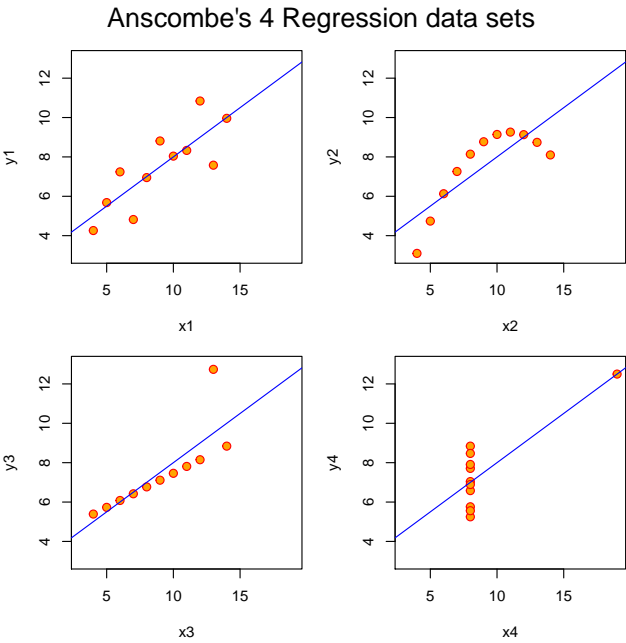


Figure 2.3: bla

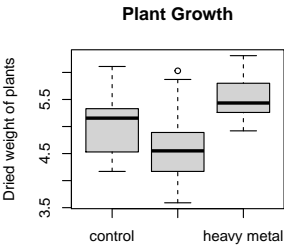


Figure 2.4: bla

details here. I recommend to look at [exercise on graphics with R](#) that accompanies this primer, as well as at

- [QuickR](#)
- [RGraphCompendium](#)
- [rexrepos](#)

Inferential statistics

Inferential statistics is concerned with inference, i.e. drawing conclusions from observations. Inference is not always, but in most cases connected to the idea with a data-generating model.

Statistical inference is drawing conclusions from observations using statistical methods

3.1 Obtaining a data-generating model

Why do we need the idea of a data-generating model to draw conclusions from data? Imagine we want to know whether plant growth can be affected by music. We might take two pots, each with a plant, and expose one to classical music and the other to heavy metal. Inevitably, one of them will grow taller, but this could be by chance, as there is always some variation in the growth rates.

Hence, we need more repetitions. Let's say we take a few more pots, 30 in the hypothetical case I show below

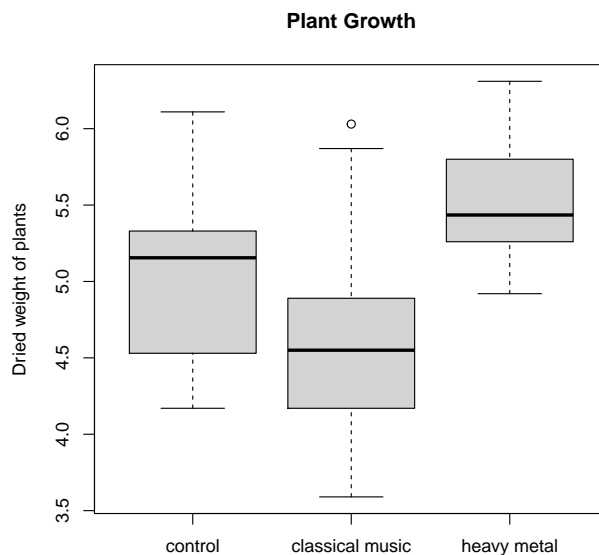


Figure 3.1: bla

There seem to be some differences between the three cases, but there is also quite a bit of variation in plant growth within each treatment (we have on average seven observations per treatment

Remember the interpretation of the boxplots - the strong line in the middle is the median. The box covers the central 0.5 quantile of the distribution.

here). Hence, it is still possible that the differences that we observe have arisen by chance.

If we want to say something definite about the probability that there is a difference between those three cases, we need to make a model that describes the stochastic variation in the data, which then allows us to calculate things such as the probability of the observed data arising by chance. These assumptions are what we call a statistical model (equivalent: stochastic process, data-generating model).

The more common class of models used for this purpose are parametric statistical models. For the data that we have here, a parametric statistical model might, for example, make the assumptions that there is a mean growth rate for each treatment (control, classical and heavy metal), but that the growth of each individual plant varies with a normal distribution around the mean growth of its respective group. The parameters of this model are the unknown mean growth rates and the variance of the normal distribution. Those parameters are then fit to the data with methods explained below, and based on the fit one can calculate, for example, the probability that the data would arise if there was no difference between the groups.

The other option to obtain a data-generating model are non-parametric methods. Non-parametric methods bypass the necessity of making assumptions, e.g. about the distribution of the data, typically by randomizing or resampling the data itself. How does this work? For the plant growth, for example, we could answer the question of how likely it is to see the observed data if there is no difference between groups also without making an assumption about the distribution - we simply throw all observations in one pot, irrespective of their treatment, and then re-distribute them randomly on the three treatments. If we do this many times (e.g. 1000 times), we can get a good idea how likely it would be to obtain the observed differences if treatment has no effect.

Non-parametric methods are an important branch of modern statistics. Their advantage is obviously that they don't make assumptions. On the other hand, parametric methods are typically much faster, and if their assumptions are correct, they are more sensitive and powerful, meaning that, with the same amount of data, they are more likely to detect an effect if it is there. For the latter two reasons, parametric methods are currently the basis of most statistical analysis.

3.2 Inferential products

Based on the data-generating model (parametric or non-parametric), we can now apply different inferential procedures to draw conclusions about our data (in our case: to decide if music makes a difference or not). In standard statistics, there are two main inferential procedures that are applied in all kinds of settings and models, and the outputs of these two procedures are: p-values and maximum likelihood estimates. A third procedure, Bayesian inference, has become

A statistical model describes how the response variable arises as a function of the predictors as well as some random (stochastic) processes

Parametric statistics uses statistical models that describe the data-generating process in terms of functions and distributions that have parameters that then need to be fit

non-parametric statistics tries to avoid making assumptions about the data-generating process. Typically, the data-generating process is emulated by using the data itself, e.g. by resampling methods.

The higher sensitivity of parametric methods relies on the fact that they make assumptions, which, in a sense, are like additional data. Of course, all results then rely on those assumptions being correct. We will talk about how to check parametric assumptions later in the chapter.

more fashionable lately. I will mention it shortly in at the end of this section

p-values

The use of p-values is connected to the inferential method of null hypothesis significance testing (NHST). The idea is the following: if we have some data observed, and we have a statistical model, we can use this statistical model to specify a fixed hypothesis about how the data did arise. For the example with the plants and music, this hypothesis could be: music has no influence on plants, all differences we see are due to random variation between individuals. Such a scenario is called the null hypothesis. Although it is very typical to use the assumption of no effect as null-hypothesis, note that it is really your choice, and you could use anything as null hypothesis, also the assumption: "classical music doubles the growth of plants". The fact that it's the analyst's choice what to fix as null hypothesis is part of the reason why there are a large number of tests available. We will see a few of them in the following chapter about important hypothesis tests.

If we have a null hypothesis, we calculate the probability that we would see the observed data or data more extreme under this scenario. This allows us to test whether our null hypothesis is compatible with the data, and thus called a hypothesis tests. We call the probability to see the data or more extreme the p-value. If the p-value falls under a certain level (the significance level α) we say the null hypothesis was rejected, and there is significant support for the alternative hypothesis. The level of α is a convention, in ecology we chose typically 0.05, so if a p-value falls below 0.05, we can reject the null hypothesis.

A problem with hypothesis tests and p-values is that their results are notoriously misinterpreted. The p-value is NOT the probability that the null hypothesis is true, or the probability that the alternative hypothesis is false, although many authors have made the mistake of interpreting it like that (?). Rather, the idea of p-values is to control the rate of false positives (Type I error). When doing hypothesis tests on random data, with an α level of 0.05, one will get exactly 5% false positives. Not more and not less.

Parameter estimates

The second type of output that is reported by most statistical methods are parameter estimates. Without going into too much detail: while the p-value is the probability of the data or more extreme data under a fixed (null) hypothesis, a parameter estimate (often also called point estimate or maximum likelihood estimate) refers to those of many possible hypothesis that are spanned by a parameter that has the highest probability to produce the observed data.

In plain words, the parameter estimate is our best estimate for a difference between groups, or the influence between parameters.

A null hypothesis H_0 is a fixed scenario that makes predictions about the expected probabilities of different observations.

The p-value is the probability to see the data or more extreme data under the null-hypothesis.

if $p < 0.05$, we say we have significant evidence to reject the null-hypothesis.

Again, in contrast, the p-value typically tells us the probability that we would see the data if there is no such influence.

Parameter estimates are typically accompanied by confidence intervals. Sloppily you can think of the 95% confidence interval of a parameter as the probable range for this area. Somewhat confusing for many, it's not the interval in which the true parameter lies with 95% probability. Rather, in a repeated experiments, the standard 95% CI will contain the true value in 95% of all cases. It's a subtle, but important difference. However, for now, don't worry about it. The CI interval is roughly where we expect the parameter.

Bayesian estimates

Finally, Bayesian methods calculate a quantity that is called the posterior parameter estimate. It is similar, but not identical to the parameter estimates discussed previously. You won't need this here, but if you want to know more about those, have a look at (?) and at my website [Learning Bayes](#).

3.3 Important hypothesis tests

After having discussed the concept of hypothesis tests (which produce p-values against a null-hypothesis) and parameter estimates (which produce best estimates and confidence intervals), let's move to practice and get to know the two probably most applied hypothesis tests, the t-test and ANOVA.

t-test

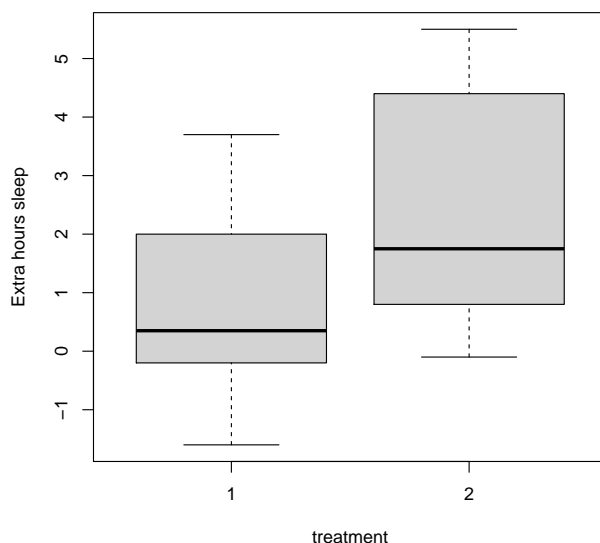
A t-test tests for differences between the means of two normally distributed samples; or if there is only one sample, between 0 and the mean of the sample. Classifying our variables, we would say the response variable is continuous, and the predictor is categorical (group 1 or group 2), or, if there is only one group, there is no predictor. Again, the statistical model underlying is that of a normally distributed response, and the null hypothesis is that there is no difference in the mean of this normal distribution for the two groups, respectively that the sample mean is 0 if we have one group. Depending on the software used, there are usually a number of adjustments possible, e.g. relaxing the assumption that the two groups have the same variance. An example in R, using the classical data from ?. The data which show the effect of two soporific drugs (increase in hours of sleep compared to control) on 10 patients.

```
> ## Traditional interface
> with(sleep, t.test(sleep$extra[sleep$group == 1], extra[group == 2]))

> ## Formula interface
> t.test(extra ~ group, data = sleep)
```

Welch Two Sample t-test

Figure 3.2: Data from ?



```
data: extra by group
t = -1.8608, df = 17.776, p-value = 0.07939
alternative hypothesis: true difference in means is not equal to 0
95 percent confidence interval:
 -3.3654832  0.2054832
sample estimates:
mean in group 1 mean in group 2
      0.75      2.33
```

Note that the output give a p-value, but also contains the estimates of the means, together with the confidence intervals.

Suggestion for reporting: $p > 0.05$: differences between groups were not significant. $p < 0.05$: we found a difference of $X \pm$ Confidence interval between the groups (p-value for difference from a t-test was X).

Analysis of variance (ANOVA)

ANOVA or analysis of variance can mean different things to different people. The standard ANOVA makes basically the same assumptions as a t-test (normally distributed responses), but allows for more than two groups. More precisely, it's a test for whether the measured response (i.e. the dependent variable) can be influence by one or several categorical variables that could have two or more levels could also interact. An interaction between two variables means that the value of one explanatory effect affects how strongly another explanatory variable affects the response. Although maybe a bit cryptic on first reading, this is pretty common. See more details [here](#).

While the word ANOVA is generally associated to the assumption of the standard ANOVA explained above (which correspond to a

linear regression, see next chapter), the concept of ANOVA can be extended in the same way as linear regression models can be extended to generalized linear models etc. Hence, we can do ANOVA also for models with non-normally distributed errors (of course you have to tell this the software, it won't do it automatically).

Here a simple example with a standard ANOVA (normal errors), testing whether weight (of chicken) depends on their diet, where diet is a factor variable with four levels:

```
> aovresult <- aov(weight~Diet, ChickWeight)
> summary(aovresult)
```

	Df	Sum Sq	Mean Sq	F value	Pr(>F)
Diet	3	155863	51954	10.81	6.43e-07 ***
Residuals	574	2758693	4806		

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

>
```

We find a p-value of 6.43e-07, which is highly significant at an α level of 0.05. Hence, we can reject the null hypothesis that the diet has no influence on the response "weight". Note that in this case, we don't get any parameter estimates, and we can't say anything about which of the diets differs from which. If you want those, there are two options:

- Either you apply what is called post-hoc testing, which means that you test for differences (e.g. with a t-test) between the diets.
- Or you switch to a regression, which is described in the next chapter

If you do post-hoc testing, you are doing multiple tests on the same data. This is a problem - the idea of the p-value is that you calculate the probability of seeing the data under ONE null hypothesis. If you do this, you will get at most 5% error at an α level of 0.05. However, if we do multiple tests, we are testing multiple null hypotheses, and there are more options for the test statistics to get significant just by chance. Hence, we need to correct the p-values for multiple testing. There are a number of options to do so, google is your friend.

When doing multiple tests on the same data, we need to correct the p-values for multiple testing.

3.4 Other important tests

t-test and ANOVA are the commonly needed tests in the context of the research skills module, but there are many more tests that could be potentially important. Important tests are

A list of tests that have wikipedia articles can be found at [here](#)

3.5 Regression models

Regression does not necessarily mean that use a different statistical model as in hypothesis testing (ANOVA and the linear regression model in R use the same assumptions). However, the goal of regression is a different one. While hypothesis tests are all about seeing whether the data would be compatible with a null-hypothesis, regression is about finding the best-fitting hypothesis or parameters. To say this again, if we have a model with a number of parameters that describe the influence of some predictors on the response, a hypothesis tests would typically set them to 0, testing the null hypothesis that there is no influence. A regression model tries to find the parameter combination that produces the highest probability to create the observed data, i.e. we are looking for the best fit.

Linear regression

The most basic regression model is the linear regression. The assumption here is that we have a response that depends on the predictors in a form of

$$y \sim a \cdot x + b + \epsilon \quad (3.1)$$

where y is the response, x is a predictor, a is the parameter that fits how strongly the predictor influences the response, b is the intercept, and ϵ is the random variation, which in a linear regression is assumed to be normally distributed.

In R, we can do such a regression by typing

```
> fit = lm(airquality$Temp~airquality$Ozone)
> summary(fit)
```

Call:

```
lm(formula = airquality$Temp ~ airquality$Ozone)
```

Residuals:

Min	1Q	Median	3Q	Max
-22.147	-4.858	1.828	4.342	12.328

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	69.41072	1.02971	67.41	<2e-16 ***
airquality\$Ozone	0.20081	0.01928	10.42	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 6.819 on 114 degrees of freedom

(37 observations deleted due to missingness)

Multiple R-squared: 0.4877, Adjusted R-squared: 0.4832

F-statistic: 108.5 on 1 and 114 DF, p-value: < 2.2e-16

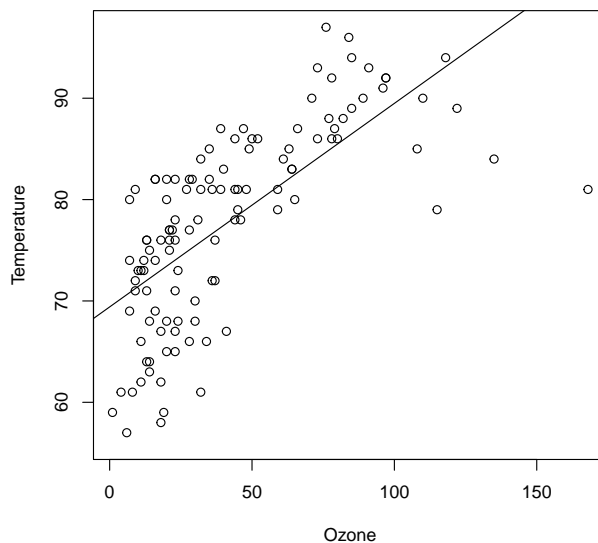


Figure 3.3: Airquality data: Temperature plotted against Ozone. Relationship indicated by the regression line.

You can use the same code regardless of whether your predictor is continuous or categorical. In case of continuous variable, a line is fit to the data. In case of a categorical variable with n levels, the first level is set as the reference (intercept), and $n-1$ parameters are fitted for the following levels that describe the difference to the reference.

The fitted parameters appear in the column "Estimate". This tells us how much the predictor, in this case Ozone, affects the response, in this case the Temperature: for each unit of Ozone more, temperature increases by a 0.208 units, with a standard error (confidence interval) of 0.019. Apart from seeing how a regression output looks like, this teaches us another valuable lesson: the fact that we have used temperature as a response here and ozone as a predictor doesn't mean that ozone causally affects temperature. In fact, it is likely the other way around: if we have more sun, it's hotter, and we tend to have more ozone as well. Regression, as most other statistical analysis, doesn't establish causality, it establishes correlation. What we are saying is that if our ozone measurements go up, we can be pretty sure that it is hotter as well. Doesn't mean that ozone creates heat. Correlation is not causality.

Correlation is not causality.

The regression results gives us a lot of p-values as well. These are results of various hypothesis tests that are performed automatically for you after the regression is done. For example, we get a p-value for each parameter. This p-value is based on a particular type of t-tests where the full model is tested against the model with the parameter set to 0. There is also other p-value, based on a different test statistics at the end of the regression output. This tests the null hypothesis that all parameters are zero.

Specifying different model assumptions in R:

Response y depends linearly on a variable a (continuous or categorical)

```
> fit = lm(y~a)
> summary(fit)
```

Response y depends linearly on two variables a and b (continuous or categorical), but the value of either variable doesn't influence the effect the other variable has on the response (no interaction)

```
> fit = lm(y~a+b)
> summary(fit)
```

Response y depends linearly on two variables a and b (continuous or categorical), but the value of one variable does influence the effect the other variable has on the response (interaction)

```
> fit = lm(y~a+b)
> summary(fit)
```

Response y depends as in $a + a^2$ on a variable a (continuous or categorical)

```
> fit = lm(y~a + I(a^2))
> summary(fit)
```

the $I()$ notation means that the following expression is interpreted as a mathematical formula.

Checking the assumptions

Formally, we can fit any data with a linear model. However, as in any statistical inference procedure the results (i.e. parameter estimates, p-values) are conditional on the assumptions that we have made. Hence, the p-value we get is conditional on the assumption that the data is actually from a process that conforms to eq. 3.2. If it doesn't the p-value could be completely wrong. Hence, we have to check whether those assumptions are actually met.

So, what were the assumptions of a linear regression? One problem I often encounter is that students remember that the assumptions were normal distribution. Hence, they look at whether the response variable is normally distributed. However, if you look sharply at eq. 3.2, you see that this was actually not the point. If we shift around the terms in eq. 3.2, we see that what is actually supposed to be normally distributed is

$$y - (a \cdot x + b) \sim \epsilon \quad (3.2)$$

i.e. the difference between the observed value and the model predictions. These differences are called the residuals, and, according to the assumptions of our model, they should be normally distributed. You get basic residual diagnostics by typing `plot(fit)`, where `fit` is your fitted model object. For further details on residual diagnostics, see [here](#).

3.6 Generalized linear regression models

The general ideas of a linear regression was that 1) The response is continuous, theoretically from -infinity to + infinity, and 2) residuals are normally distributed around the model predictions. The idea of the GLM framework is take the linear regression framework is to allow you to work as before in the linear regression example, but relaxing both the assumptions about response values from - to + infinity, and the normality. To do this, we have to do two things

- To get the output values on the range that we want, we wrap the linear model in a transformation function that forces the response on the right interval (typical intervals are positive, or between 0 and 1). This transformation is called the link function
- To fit other distributions, we have to tell the model to use something else than the Gaussian error function.

Let's talk about these points in a bit more detail.

The link function

We said above that a linear regression takes the form

$$y \sim a \cdot x + b \quad (3.3)$$

That means that if x gets large, y could take any value, positive or negative. A trick to ensure that all predictions for y are positive, or within a certain range is using a link function of the form

$$y \sim f^{link}(a \cdot x + b) \quad (3.4)$$

Any function is possible, but as we see later typical choices are the exponential function, which ensures positives outcomes, and the inverse logit, which ensures are range between 0 and 1.

Other distributions

Well, this is conceptually the easy part, but maybe you are not yet aware what kind of distributions exist beside the normal. Two typical choices that we use below are the binomial (the distribution for coin flipping), and the Poisson distribution (a discrete probability distribution). There are many other choices available. Maybe it becomes more clear when we move to the actual examples in the next sections.

0/1 data - logistic regression

Logistic regression is the most common analysis for binary data (presence/absence; survived/dead; infected/not infected). Logistic regression assumes that the distribution is binomial (coin flip model). To get the linear predictor on a scale between 0 and 1 that is necessary for the binomial distribution, we use the logistic link function (or inverse logit).

evspacelcm

In R:

Here an example with the data of the Titanic survivors. Note that if you tell R to use the binomial distribution, the logit link is automatically selected. If you wanted, you could overrule this choice.

```
> library(effects)
> fmt <- glm(survived ~ age + I(age^2) + I(age^3), family=binomial, data = TitanicSurvival)
> summary(fmt)
```

Call:

```
glm(formula = survived ~ age + I(age^2) + I(age^3), family = binomial,
     data = TitanicSurvival)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-1.5062	-0.9978	-0.9695	1.3483	2.0135

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	7.855e-01	3.031e-01	2.592	0.009549	**
age	-1.189e-01	3.291e-02	-3.613	0.000303	***
I(age^2)	3.414e-03	1.113e-03	3.066	0.002171	**
I(age^3)	-2.931e-05	1.107e-05	-2.648	0.008109	**

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 1414.6 on 1045 degrees of freedom

Residual deviance: 1398.5 on 1042 degrees of freedom

(263 observations deleted due to missingness)

AIC: 1406.5

Number of Fisher Scoring iterations: 4

count data - poisson regression

Poisson regression is the standard choice for working with count data, although there are a few other options available as well. In poisson regression, the standard choice is to use an exponential function to make all values positive. The inverse of the exponential is the log, so we call this the log link. As before, R is choosing this automatically if you specify the distribution to be poisson.

In R:

An example, using some data on the feeding of bird nestlings, in relation to their attractiveness:

```
> schnaepper <- read.csv("schnaepper.txt", sep="")
> fm <- glm(stuecke ~ attrakt, family=poisson, data = schnaepper)
```

```

> summary(fm)

Call:
glm(formula = stuecke ~ attrakt, family = poisson, data = schnaepper)

Deviance Residuals:
    Min       1Q   Median       3Q      Max
-1.55377  -0.72834   0.03699   0.59093   1.54584

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  1.47459     0.19443   7.584 3.34e-14 ***
attrakt      0.14794     0.05437   2.721  0.00651 **
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for poisson family taken to be 1)

    Null deviance: 25.829  on 24  degrees of freedom
Residual deviance: 18.320  on 23  degrees of freedom
AIC: 115.42

Number of Fisher Scoring iterations: 4

```

Residual checks in GLMs

Residuals in glms are not supposed to be normally distributed, so don't use standard checks for normality such as normal qq plots to check for the appropriateness of the residuals. For not too complicated models, a good way to deal with this problem is to use the so-called pearsons residuals, which scale the observed differences between model and data by the expected variance of the model ¹

One standard concern in poisson or binomial glms is that the variance of the poisson and binomial distribution cannot be adjusted, but is fixed by the mean. This is a problem you don't encounter in the normal linear model, because here the random part is modelled by a normal distribution, which has a parameter for the variance. A problem that appears very frequently in Poisson or binomial glms is overdispersion, i.e. that the residuals show more variance than expected under the fitted model. The easiest way to correct for this is to use the quasi-Poisson and quasi-binomial models available in the glm function. These models fit an additional parameter that modifies the variance of the Poisson and binomial glm.

¹ in R, you can specify the option `pearson` in many functions, including the `residual()` function that you can apply to a fitted object

You can check for overdispersion by looking at the fitted deviance, or apply an overdispersion test

Predictive statistics - machine learning

A third class of statistical procedures that have become very important in recent years are predictive methods, often called machine-learning algorithms. The basic goal of these methods is to be able to make predictions from a given dataset with the lowest possible error. In doing so, they typically use relatively complicated, often non-parametric methods that typically don't allow to calculate inferential products such as the MLE or p-values.

There is considerable tension between the more classical field of inferential statistics and the more modern field of machine-learning. For classical, inferential statisticians, machine learning methods have abandoned the idea of "learning from data", in the sense of comparing hypotheses to data, in favor of simply making predictions. A statistician concentrating on machine-learning would reply that in many applied problems, there is nothing to learn¹. The goal is to build an algorithm that is able to correctly predict given a complex dataset. The distinction between the goals of inferential statistics and predictive statistics, as well as the tension between these fields, are nicely summarized in the abstract of the extremely recommendable article "Statistical Modeling: The Two Cultures" by ?:

¹ Typical machine learning applications include predicting the interests of customers in web shops, the association of feature-rich satellite data with ground signals, or speech / face recognition. Machine-learning experts are currently sought-after by technology companies such as google, facebook and so on.

There are two cultures in the use of statistical modeling to reach conclusions from data. One assumes that the data are generated by a given stochastic data model. The other uses algorithmic models and treats the data mechanism as unknown. The statistical community has been committed to the almost exclusive use of data models. This commitment has led to irrelevant theory, questionable conclusions, and has kept statisticians from working on a large range of interesting current problems. Algorithmic modeling, both in theory and practice, has developed rapidly in fields outside statistics. It can be used both on large complex data sets and as a more accurate and informative alternative to data modeling on smaller data sets. If our goal as a field is to use data to solve problems, then we need to move away from exclusive dependence on data models and adopt a more diverse set of tools.

I have included this short chapter because of the importance of predictive methods in modern statistics. A detailed explanation of the methods of machine-learning, however, is beyond this introduction. If you are interested in starting to learn more about predictive

methods, I would recommend to start with the textbook by ? that I also recommend at the end of this primer for further reading.

5

Design of experiments

Let's come back to one of the first point in this script: the data. If we have to collect data ourselves, we have to answer a number of questions. Which variables should we collect? At which values of those variables should we collect data? And how many replicates do we need?

5.1 Which variables

In a practical setting, we are typically interested in how a response is affected by a number of predictor variables. Clearly, we need to measure both response and this predictors of interest across a few of those predictor values to say something about the effect of the predictors. If we only wanted to know whether there is a correlation between predictors and response, our list of variables would be complete at this point. However, typically, we want to know not only if there is a correlation, but also whether we can say with some confidence that this correlation is causal. If we want to make this claim, we have to exclude that there are founfounding factors, also called confounding variables.

Correlation is not causality. For suggesting causality, we have to exclude confounding effects.

What is a confounding variable?

Imagine we are interested in a response A, and we have hypothesized that A B. Imagine there is a second predictor variable C that has an influence on A, but in which we are not interested in for the purpose of the question under consideration. Such a variable that is not of interest for the question is also called "extraneous variables". . So we also have A C, but we are not interested in this relationship. If we now take data, and don't measure C, it's usually not a bit problem as long as C is uncorrelated with B - it might create a bit more variability in the response, but by and large the effect of C should average out and we should still be able to detect the effect of B.

An extraneous variable is a variable that can influence the response, but is not of interest for the experimenter

The problem of confounding appears when the extraneous variable C is for some reason correlated with the predictor variable of interest B. In that case, if we only measure B, we see both the effect of B and C. In this case, we may attribute the effect of C on A wrongly to the effect of B on A. Auch a correlation that is caused by an

A confounding variable is an extraneous variable that correlated to both the response and a predictor variable of interest.

A spurious correlation is a correlation that is caused by a confounding variable.

unmeasured confounding variable is called a spurious correlation.

What do do about confounding variables

If we think there is a factor that could be confounding, we basically have three options

1. Best: control the value of these factors. Either fix the value (preferred if we are not interested in this factor), else vary the value in a controlled way.
2. Second best: randomize and measure them
3. Third best: only randomize or only measure them

Randomization means that we try to ensure that the confounding factor is not systematically correlated with the variable of interest (but can still cause problems with interactions and nonlinear relationships).

Measuring allows to control the effect, but cost power (see below) and, and we can't measure everything

5.2 At which values should we measure the variables

If we have decided which variables to measure, we have to decide for the values at which we want to measure them. In an observational study, this may not always be possible completely, but it's usually possible to ensure sufficient variation in the predictor variables. A few points to consider

Vary all variables independently

A common problem in practice is that we have two variables, but their values change in a correlated way. Imagine we test for the presence of a species, but we have only warm dry and cold wet sites. We say the two variables are collinear. In this case we don't know whether any observed effect is due to temperature or water availability. The bottom-line: if you want to separate two effects, you need them to vary independently.

Interactions

To be able to detect interactions between variables, it's not enough to vary all, you also need to have certain combinations. The buzzword here is called factorial design. Google will help you.

Nonlinear effects

The connection of two points is a line. If you want to see whether the response to a variable is nonlinear, you therefore need more than two values of each variable.

5.3 How many replicates?

We said before that the significance level α is the probability of finding false positives. This is called the type I error. There is another error we can make: failing to find significance for a true effect. This is called the type II error, and the probability of finding an effect is called power. . For standard statistical methods, power can be calculated. You have to look it up for your particular method, but in general assume that

Power is the probability of finding significance for an effect if it's there

1. Power goes up with increasing effect size
2. Power goes down with increasing variability in the response

This means that, unlike for the type I error which is fixed, calculation of power requires knowledge about the expected effect and the variability. This sounds really bad, but in most cases you can estimate from previous experience how much variation there will be, and in most cases you also know how big the effect has to be at least to be interesting. Based on that, you can then calculate how many samples you need.

Checklist experimental design

- () Clear, logically consistent question? Write it down.
- () What variables need to be measured to answer this question? Check that the measured variables really correspond to the main question
- () Confounding factors controlled, randomised or measured? Are you sure they are confounding (correlated to response AND one or several of the predictors)
- () Define the exact hypothesis (statistical model) to be tested. Write it down, as in $height \sim age + soil * precipitation + precipitation^2$. Based on that, decide how many different levels of each variable need to be measured.
- () Decide on the number of replicates. Make a guess for effect size and variability of the data, and either calculate or guess the number of replicates necessary to get sufficient power. What sufficient means depends on the field, but I would say you want to have a good chance to see an effect if it's there, so a power of $> 80\%$ would be good.

6

Good to know and further reading

6.1 Reproducibility and good scientific practice

Reproducibility means that each step of your analysis is repeatable. Experience shows that it is not as trivial as it sounds to ensure reproducibility. Here some hints for making your data analysis reproducible

- Once you have your raw data produced, NEVER change it. Store it in a safe location, make a backup, and never touch it again
- Typically you will have to do some cleaning, renaming etc. before the data analysis. If possible at all, make this through a script (e.g. R, python, perl). Store the script with the analysis.
- Use a version control system for your code, and note for each output the revision number that the output was produced with.
- When running the analysis, store the random seed and the settings of your computer to ensure reproducibility. In R, the easiest way to do this is to set the random seed by `random.seed(123)`, and store the results of `sessionInfo()` which provides you with the version numbers of all the packages that you use
- Think about running your code within an reporting environment such as Rmd or sweave

6.2 How to learn more about statistics

- To complete this primer, I would recommend that you go through the [practicals for this script](#).
- If you want one further practical textbook for beginners, I recommend ? for German speakers (ebook free of charge for students from Freiburg, contact me) and ? for English speakers.
- For the technically slightly more ambitious (it's still very elementary), I recommend ?. You can download the pdf for free, and there is a MOOC available for the book with lectures and exercises.

- For more help and references, see the stats help website of our department [here](#), in particular the recommendations regarding R scripts and statistics textbooks.

7

Bibliography