

Statistical Modelling

Contents

Preliminary remarks	5
1 Introduction to statistical inference	7
1.1 Hypothesis testing	8
1.2 Exploratory Data Analysis	20
2 Linear regression	35
2.1 Introduction	36
2.2 Ordinary least squares	38
2.3 Interpretation of the model parameters	40
2.4 Tests for parameters of the linear model	47
2.5 Coefficient of determination	52
2.6 Predictions	53
2.7 Interactions	54
2.8 Collinearity	64
2.9 Graphical analysis of residuals	67
3 Likelihood-based inference	83
3.1 Maximum likelihood	83
3.2 Likelihood-based tests	88
3.3 Profile likelihood	91
3.4 Information criteria	93
4 Generalized linear models	97
4.1 Basic principles	97
4.2 Theory of generalized linear models	98
5 Correlated and longitudinal data	103
6 Linear mixed models	105

7	Survival analysis	107
A	Additional topics and prerequisites	109
A.1	Population and samples	109
A.2	Random variable	110
A.3	Laws of large numbers	118
A.4	Central Limit Theorem	119
B	Mathematical derivations	123
B.1	Derivation of the ordinary least squares estimator	123
B.2	Derivation of the coefficient of determination	124
C	R	125
C.1	Basics of R	125
C.2	Linear models in R using the <code>lm</code> function	129

Preliminary remarks

These notes by Léo Belzile (HEC Montréal) are licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International License and were last compiled on 2020-10-04.

While we show how to implement statistical tests and models in **SAS** in class, these note will illustrate the concepts using **R**: visit the R-project website to download the program. The most popular graphical cross-platform front-end is RStudio Desktop.

The most famous quote about statistical models is probably due to George Box, who claimed that “all models are wrong, but some are useful”. This standpoint is reductive: Peter McCullagh and John Nelder wrote in the preamble of their book (emphasis mine)

Modelling in science remains, partly at least, an art. Some principles do exist, however, to guide the modeller. The first is that all models are wrong; **some, though, are better** than others and we can **search for the better ones**. At the same time we must recognize that eternal truth is not within our grasp.

And this quote by David R. Cox adds to the point:

... it does not seem helpful just to say that all models are wrong. The very word model implies simplification and idealization. The idea that complex physical, biological or sociological systems can be exactly described by a few formulae is patently absurd. The construction of idealized representations that **capture important stable aspects of such systems** is, however, a vital part of general scientific analysis and statistical models, especially substantive ones, do not seem essentially different from other kinds of model.

Chapter 1

Introduction to statistical inference

Statistical modelling requires a good grasp of statistical inference: as such, we begin with a review of hypothesis testing and graphical exploratory data analysis.

The purpose of statistical inference is to draw conclusions based on data. Scientific research relies on hypothesis testing: once an hypothesis is formulated, the researcher collects data, performs a test and concludes as to whether there is evidence for the proposed theory.

There are two main data type: **experimental** data are typically collected in a control environment following a research protocol with a particular experimental design: they serve to answer questions specified ahead of time. This approach is highly desirable to avoid the garden of forking paths (researchers unfortunately tend to refine or change their hypothesis in light of data, which invalidates their findings — preregistration alleviates this somewhat). While experimental data are highly desirable, it is not always possible to collect experimental data: for example, an economist cannot modify interest rates to see how it impacts consumer savings. When data have been collected beforehand without intervention (for other purposes), these are called **observational**. These will be the ones most frequently encountered.

A stochastic model will comprise two ingredients: a distribution for the random data and a formula linking the parameters or the conditional expectation of a response variable Y to a set of explanatories X . A model can serve to either predict new outcomes (predictive modelling) or else to test research hypothesis about the effect of the explanatory variables on the response (explanatory model). These two objectives are of course not mutually exclusive even if we distinguish in practice inference and prediction.

A predictive model gives predictions of Y for different combinations of explanatory variables or future data. For example, one could try to forecast the energy consumption of a house as a function of weather, the number of inhabitants and its size. Black boxes used in machine

learning are often used solely for prediction: these models are not easily interpreted and they often ignore the data structure.

By contrast, explicative models are often simple and interpretable: regression models are often used for inference purpose and we will focus on these.

- Are consumer ready to spend more when they pay by credit card rather than by cash?
- Is there wage discrimination towards women in a US college?
- University degree: “is the university experience worth the cost’ ”?
- What are the criteria impacting health insurance premiums?
- Is the price of gasoline more expensive in the Gaspé peninsula than in the rest of Quebec? A report of the *Régie de l’énergie* examines the question
- Are driving tests in the UK easier if you live in a rural area? An analysis of *The Guardian* hints that it is the case.
- Does the risk of transmission of Covid19 increase with distancing? A (bad) meta-analysis says two meters is better than one (or how to draw erroneous conclusions from a bad model).

1.1 Hypothesis testing

An hypothesis test is a binary decision rule used to evaluate the statistical evidence provided by a sample to make a decision regarding the underlying population. The main steps involved are:

- define the model parameters
- formulate the alternative and null hypothesis
- choose and calculate the test statistic
- obtain the null distribution describing the behaviour of the test statistic under \mathcal{H}_0
- calculate the p -value
- conclude (reject or fail to reject \mathcal{H}_0) in the context of the problem.

A good analogy for hypothesis tests is a trial for murder on which you are appointed juror.

- The judge lets you choose between two mutually exclusive outcome, guilty or not guilty, based on the evidence presented in court.
- The presumption of innocence applies and evidences are judged under this optic: are evidence remotely plausible if the person was innocent? The burden of the proof lies with the prosecution to avoid as much as possible judicial errors. The null hypothesis \mathcal{H}_0 is *not guilty*, whereas the alternative \mathcal{H}_a is *guilty*. If there is a reasonable doubt, the verdict of the trial will be not guilty.
- The test statistic (and the choice of test) represents the summary of the proof. The more overwhelming the evidence, the higher the chance the accused will be declared guilty. The prosecutor chooses the proof so as to best outline this: the choice of

evidence (statistic) ultimately will maximise the evidence, which parallels the power of the test.

- The final step is the verdict. This is a binary decision, guilty or not guilty. For an hypothesis test performed at level α , one would reject (guilty) if the p -value is less than α .

The above description provides some heuristic, but lack crucial details developed in the next section written by Juliana Schulz.

1.1.1 Hypothesis

In statistical tests we have two hypotheses: the null hypothesis (H_0) and the alternative hypothesis (H_1). Usually, the null hypothesis is the ‘status quo’ and the alternative is what we’re really interested in testing. A statistical hypothesis test allows us to decide whether or not our data provides enough evidence to reject H_0 in favour of H_1 , subject to some pre-specified risk of error. Usually, hypothesis tests involve a parameter, say θ , which characterizes the underlying distribution at the population level and whose value is unknown. A two-sided hypothesis test regarding a parameter θ has the form

$$\mathcal{H}_0 : \theta = \theta_0 \quad \text{versus} \quad \mathcal{H}_a : \theta \neq \theta_0.$$

We are testing whether or not θ is precisely equal to the value θ_0 . The hypotheses are a statistical representation of our research question.

For example, for a two-sided test for the regression coefficient β_j associated to an explanatory variable X_j , the null and alternative hypothesis are explicative d’intérêt X_j , les hypothèses sont

$$\mathcal{H}_0 : \beta_j = \beta_j^0 \quad \text{versus} \quad \mathcal{H}_a : \beta_j \neq \beta_j^0,$$

where β_j^0 is some value that reflects the research question of interest. For example, if $\beta_j^0 = 0$, the underlying question is: is covariate X_j impacting the response Y once other variables have been taken into account?

Note that we can impose direction in the hypotheses and consider alternatives of the form $\mathcal{H}_a : \theta > \theta_0$ or $\mathcal{H}_a : \theta < \theta_0$.

1.1.2 Test statistic

A test statistic T is a functional of the data that summarise the information contained in the sample for θ . The form of the test statistic is chosen such that we know its underlying distribution under H_0 , that is, the potential values taken by T and their relative probability if H_0 is true. Indeed, Y is a random variable and its value change from one sample to the

next. This allows us to determine what values of T are likely if H_0 is true. Many statistics we will consider are **Wald statistic**, of the form

$$T = \frac{\hat{\theta} - \theta_0}{\text{se}(\hat{\theta})}$$

where $\hat{\theta}$ is an estimator of θ , θ_0 is the postulated value of the parameter and $\text{se}(\hat{\theta})$ is an estimator of the standard deviation of the test statistic $\hat{\theta}$.

For example, to test whether the mean of a population is zero, we set

$$\mathcal{H}_0 : \mu = 0, \quad \mathcal{H}_a : \mu \neq 0,$$

and the Wald statistic is

$$T = \frac{\bar{X} - 0}{S_n/\sqrt{n}}$$

where \bar{X} is the sample mean of X_1, \dots, X_n ,

$$\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i = \frac{X_1 + \dots + X_n}{n}$$

and the standard error (of the mean) \bar{X} is S_n/\sqrt{n} ; the sample variance S_n is an estimator of the standard deviation σ ,

$$S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2.$$

It's important to distinguish between procedures/formulas and their numerical values. An **estimator** is a rule or formula used to calculate an estimate of some parameter or quantity of interest based on observed data. For example, the sample mean \bar{X} is an estimator of the population mean μ . Once we have observed data we can actually compute the sample mean, that is, we have an estimate — an actual value. In other words,

- an estimator is the procedure or formula telling us how to use sample data to compute an estimate. It's a random variable since it depends on the sample.
- an estimate is the numerical value obtained once we apply the formula to observed data

1.1.3 Null distribution and p -value

The p -value allows us to decide whether the observed value of the test statistic T is plausible under H_0 . Specifically, the p -value is the probability that the test statistic is equal or more

extreme to the estimate computed from the data, assuming H_0 is true. Suppose that based on a random sample X_1, \dots, X_n we obtain a statistic whose value $T = t$. For a two-sided test $\mathcal{H}_0 : \theta = \theta_0$ vs. $\mathcal{H}_a : \theta \neq \theta_0$, the p -value is $\Pr_0(|T| \geq |t|)$. If the distribution of T is symmetric around zero, the p -value is

$$p = 2 \times \Pr_0(T \geq |t|).$$

Consider the example of a two-sided test involving the population mean $H_0 : \mu = 0$ against the alternative $H_1 : \mu \neq 0$. Assuming the random sample comes from a normal (population) $\text{No}(\mu, \sigma^2)$, it can be shown that if H_0 is true (that is, if $\mu = 0$), the test statistic

$$T = \frac{\bar{X}}{S/\sqrt{n}}$$

follows a Student- t distribution with $n - 1$ degrees of freedom, denoted St_{n-1} . This allows us to calculate the p -value (either from a table, or using some statistical software). The Student- t distribution is symmetric about zero, so the p -value is $P = 2 \times \Pr(T_{n-1} > |t|)$, where $T \sim \text{St}_{n-1}$.

1.1.4 Conclusion

The p -value allows us to make a decision about the null hypothesis. If \mathcal{H}_0 is true, the p -value follows a uniform distribution. Thus, if the p -value is small, this means observing an outcome more extreme than $T = t$ is unlikely, and so we're inclined to think that H_0 is not true. There's always some underlying risk that we're making a mistake when we make a decision. In statistic, there are two type of errors:

- type I error: we reject H_0 when H_0 is true,
- type II error: we fail to reject H_0 when H_0 is false.

These hypothesis are not judged equally: we seek to avoid error of type I (judicial errors, corresponding to condemning an innocent). To prevent this, we fix a the level of the test, α , which captures our tolerance to the risk of committing a type I error: the higher the level of the test α , the more often we will reject the null hypothesis when the latter is true. The value of $\alpha \in (0, 1)$ is the probability of rejecting \mathcal{H}_0 when \mathcal{H}_0 is in fact true,

$$\alpha = \Pr_0(\text{reject } \mathcal{H}_0).$$

The level α is fixed beforehand, typically 1%, 5% or 10%. Keep in mind that the probability of type I error is α only if the null model for \mathcal{H}_0 is correct (sic) and correspond to the data generating mechanism.

The focus on type I error is best understood by thinking about medical trial: you need to prove a new cure is better than existing alternatives drugs or placebo, to avoid extra

costs or harming patients (think of Didier Raoult and his unsubstantiated claims that hydrochloroquine, an antipaludean drug, should be recommended treatment against Covid19).

Decision \ true model	\mathcal{H}_0	\mathcal{H}_a
fail to reject \mathcal{H}_0	✓	type II error
reject \mathcal{H}_0	type I error	✓

To make a decision, we compare our p -value P with the level of the test α :

- if $P < \alpha$, we reject \mathcal{H}_0 ;
- if $P \geq \alpha$, we fail to reject \mathcal{H}_0 .

Do not mix up level of the test (probability fixed beforehand by the researcher) and the p -value. If you do a test at level 5%, the probability of type I error is by definition α and does not depend on the p -value. The latter is conditional probability of observing a more extreme likelihood given the null distribution \mathcal{H}_0 is true.

1.1.5 Power

There are two sides to an hypothesis test: either we want to show it is not unreasonable to assume the null hypothesis, or else we want to show beyond reasonable doubt that a difference or effect is significative: for example, one could wish to demonstrate that a new website design (alternative hypothesis) leads to a significant increase in sales relative to the status quo. Our ability to detect these improvements and make discoveries depends on the power of the test: the larger the power, the greater our ability to reject \mathcal{H}_0 when the latter is false.

Failing to reject \mathcal{H}_0 when \mathcal{H}_a is true corresponds to the definition of type II error, the probability of which is $1 - \gamma$, say. The **power of a test** is the probability of rejecting \mathcal{H}_0 when \mathcal{H}_0 is false, i.e.,

$$\gamma = \Pr_a(\text{reject } \mathcal{H}_0)$$

Depending on the alternative models, it is more or less easy to detect that the null hypothesis is false and reject in favor of an alternative.

We want a test to have high power, i.e., that γ be as close to 1 as possible. Minimally, the power of the test should be α because we reject the null hypothesis α fraction of the time even when \mathcal{H}_0 is true. Power depends on many criteria, notably

- the effect size: the bigger the difference between the postulated value for θ_0 under \mathcal{H}_0 and the observed behavior, the easier it is to detect it. (Figure 1.3);

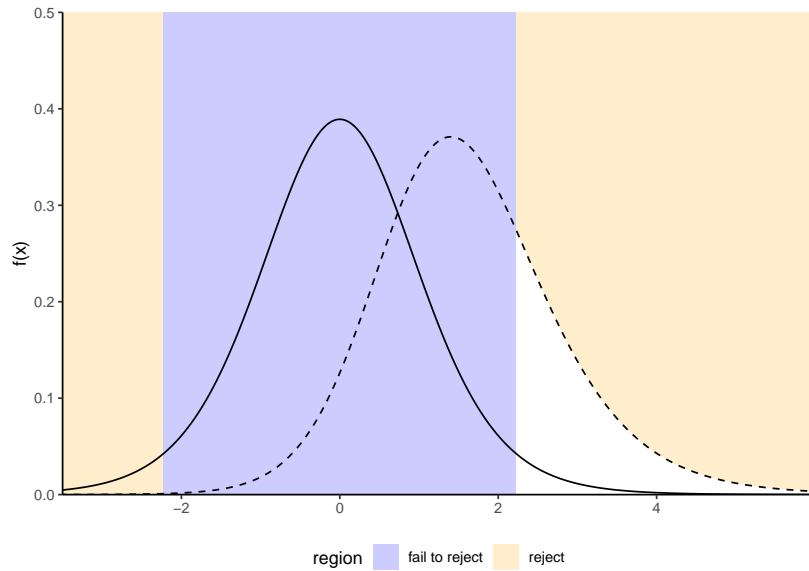


Figure 1.1: Comparison between null distribution (full curve) and a specific alternative for a t -test (dashed line). The power corresponds to the area under the curve of the density of the alternative distribution which is in the rejection area (in white).

- variability: the less noisy your data, the easier it is to detect differences between the curves (big differences are easier to spot, as Figure 1.2 shows);
- the sample size: the more observation, the higher our ability to detect significant differences because the standard error decreases with sample size n at a rate (typically) of $n^{-1/2}$. The null distribution also becomes more concentrated as the sample size increases.
- the choice of test statistic: for example, rank-based statistics discard information about the actual values and care only about relative ranking. Resulting tests are less powerful, but are typically more robust to model misspecification and outliers. The statistics we will choose are standard and amongst the most powerful: as such, we won't dwell on this factor.

To calculate the power of a test, we need to single out a specific alternative hypothesis. In very special case, analytic derivations are possible: for example, the one-sample t -test statistic $T = \sqrt{n}(\bar{X}_n - \mu_0)/S_n \sim \mathcal{T}_{n-1}$ for a normal sample follows a noncentral Student- t distribution with noncentrality parameter Δ if the expectation of the population is $\Delta + \mu_0$. In general, such closed-form expressions are not easily obtained and we compute instead the power of a test through Monte Carlo methods. For a given alternative, we simulate repeatedly samples from the model, compute the test statistic on these new samples and

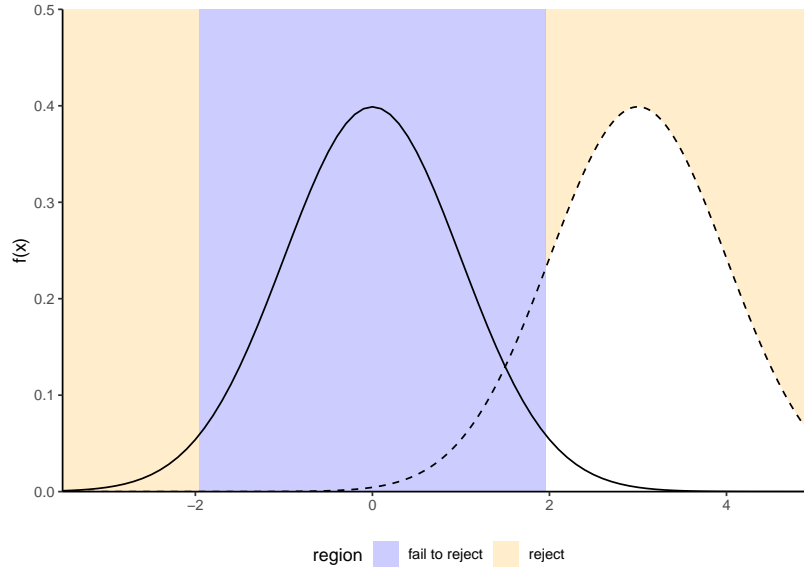


Figure 1.2: Increase in power due to an increase in the mean difference between the null and alternative hypothesis. Power is the area in the rejection region (in white) under the alternative distribution (dashed): the latter is more shifted to the right relative to the null distribution (full line).

the associated p -values based on the postulated null hypothesis. We can then calculate the proportion of tests that lead to a rejection of the null hypothesis at level α , namely the percentage of p -values smaller than α .

1.1.6 Confidence interval

A **confidence interval** is an alternative way to present the conclusions of an hypothesis test performed at significance level α . It is often combined with a point estimator $\hat{\theta}$ to give an indication of the variability of the estimation procedure. Wald-based $(1 - \alpha)$ confidence intervals for a parameter θ are of the form

$$\hat{\theta} \pm q_{\alpha/2} \text{se}(\hat{\theta})$$

where $q_{\alpha/2}$ is the $1 - \alpha/2$ quantile of the null distribution of the Wald statistic

$$T = \frac{\hat{\theta} - \theta}{\text{se}(\hat{\theta})},$$

and where θ represents the postulated value for the fixed, but unknown value of the parameter. The bounds of the confidence intervals are random variables, since both $\hat{\theta}$ and $\text{se}(\hat{\theta})$

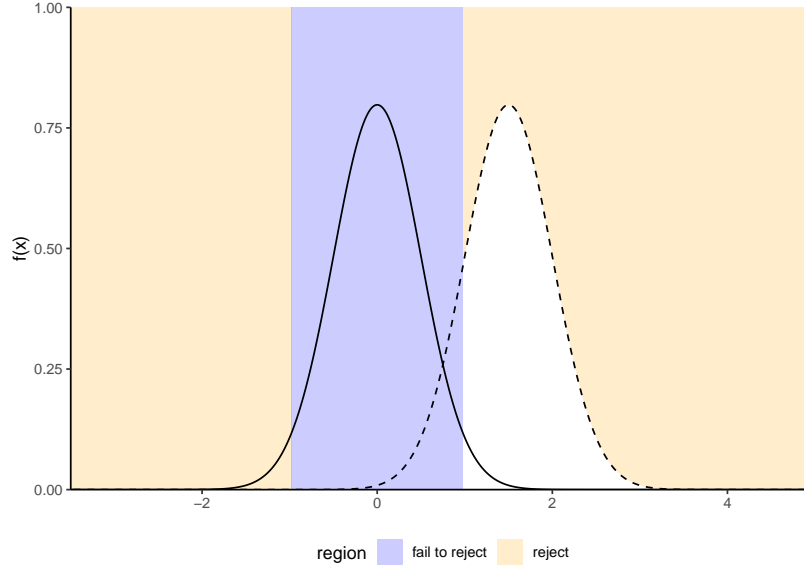


Figure 1.3: Increase of power due to an increase in the sample size or a decrease of standard deviation of the population: the null distribution (full line) is more concentrated. Power is given by the area (white) under the curve of the alternative distribution (dashed). In general, the null distribution changes with the sample size.

are random variables: their values depend on the sample, and will vary from one sample to another.

For example, for a random sample X_1, \dots, X_n from a normal distribution $\text{No}(\mu, \sigma)$, the $(1 - \alpha)$ confidence interval for the population mean μ is

$$\bar{X} \pm t_{n-1, \alpha/2} \frac{S}{\sqrt{n}}$$

where $t_{n-1, \alpha/2}$ is the $1 - \alpha/2$ quantile of a Student- t distribution with $n - 1$ degrees of freedom.

Before the interval is calculated, there is a $1 - \alpha$ probability that θ is contained in the **random** interval $(\hat{\theta} - q_{\alpha/2} \text{se}(\hat{\theta}), \hat{\theta} + q_{\alpha/2} \text{se}(\hat{\theta}))$, where $\hat{\theta}$ denotes the estimator. Once we obtain a sample and calculate the confidence interval, there is no more notion of probability: the true value of the parameter θ is either in the confidence interval or not. We can interpret confidence interval's as follows: if we were to repeat the experiment multiple times, and calculate a $1 - \alpha$ confidence interval each time, then roughly $1 - \alpha$ of the calculated confidence intervals would contain the true value of θ in repeated samples (in the same way, if you flip a coin, there is roughly a 50-50 chance of getting heads or tails, but any outcome

will be either). Our confidence is in the *procedure* we use to calculate confidence intervals and not in the actual values we obtain from a sample.

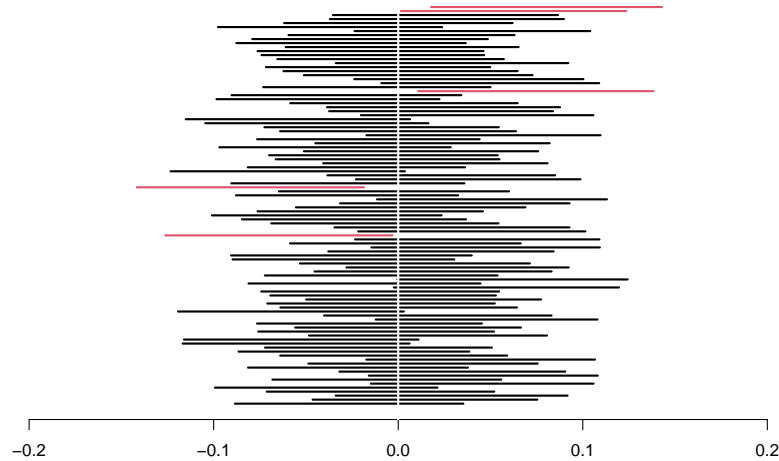


Figure 1.4: 95% confidence intervals for the mean of a standard normal population $N(0, 1)$, with 100 random samples. On average, 5% of these intervals fail to include the true mean value of zero (in red).

If we are only interested in the binary decision rule reject/fail to reject \mathcal{H}_0 , the confidence interval is equivalent to a p -value since it leads to the same conclusion. Whereas the $1 - \alpha$ confidence interval gives the set of all values for which the test statistic doesn't provide enough evidence to reject \mathcal{H}_0 at level α , the p -value gives the probability under the null of obtaining a result more extreme than the postulated value and so is more precise for this particular value. If the p -value is smaller than α , our null value θ will be outside of the confidence interval and vice-versa.

Example 1.1 (Online purchases of millenials). Suppose a researcher studies the evolution of online sales in Canada. She postulates that generation Y members make more online purchase than older generations. A survey is sent to a simple random sample of $n = 500$ individuals from the population with 160 members of generation Y and 340 older people. The response ariable is the total amount of online goods purchased in the previous month (in dollars).

In this example, we consider the difference between the average amount spent by Y members and those of previous generations: the mean difference in the samples is -16.49 dollars and thus millenials spend more. However, this in itself is not enough to conclude that the different is significant, nor can we say it is meaningful. The amount spent online varies from one individual to the next (and plausibly from month to month), and so different

random samples would yield different mean differences.

The first step of our analysis is defining the parameters corresponding to quantities of interest and formulating the null and alternative hypothesis as a function of these parameters. We will consider a test for the difference in mean of the two populations, say μ_1 for the expected amount spent by generation Y and μ_2 for older generations, with respective standard errors σ_1 and σ_2 . We next write down our hypothesis: the researcher is interested in whether millenials spend more, so this is the alternative hypothesis, $\mathcal{H}_a : \mu_1 > \mu_2$. The null consists of all other values $\mathcal{H}_0 : \mu_1 \leq \mu_2$, but only $\mu_1 = \mu_2$ matters for the purpose of testing (why?)

The second step is the choice of test statistic. We consider the Welch (1947) statistic for a difference in mean between two samples,

$$T = \frac{\bar{X}_1 - \bar{X}_2}{\left(\frac{S_1^2}{n_1} + \frac{S_2^2}{n_2}\right)^{1/2}},$$

where \bar{X}_i is the sample mean, S_i^2 is the unbiased variance estimator and n_i is the sample size for group i ($i = 1, 2$). If the mean difference between the two samples is zero, then $\bar{X}_1 - \bar{X}_2$ has mean zero and the difference has variance $\sigma_1^2/n_1 + \sigma_2^2/n_2$. For our sample, the value of statistic is $T = -2.76$. Since the value changes from one sample to the next, we need to determine if this value is compatible with the null hypothesis by comparing it to the null distribution of T (when \mathcal{H}_0 is true and $\mu_1 - \mu_2 = 0$). We perform the test at level $\alpha = 0.05$.

The third step consists in obtaining a benchmark to determine if our result is extreme or unusual. To make comparisons easier, we standardize the statistic so its has mean zero and variance one under the null hypothesis $\mu_1 = \mu_2$, so as to obtain a dimensionless measure whose behaviour we know for large sample. The (mathematical) derivation of the null distribution is beyond the scope of this course, and will be given in all cases. Asymptotically, T follows a standard normal distribution $\text{No}(0, 1)$, but there exists a better finite-sample approximation when n_1 or n_2 is small; we use Satterthwaite (1946) and a Student- t distribution as null distribution.

It only remains to compute the p -value. If the null distribution is well-specified and \mathcal{H}_0 is true, then the random variable P is uniform on $[0, 1]$; we thus expect to obtain under the null something larger than 0.95 only 5% of the time for our one-sided alternative since we consider under \mathcal{H}_0 the event $\Pr(T > t)$. The p -value is 1 and, at level 5%, we reject the null hypothesis to conclude that millenials spend significantly than previous generation for monthly online purchases, with an estimated average difference of -16.49.

Example 1.2 (Price of Spanish high speed train tickets). The Spanish national railway company, Renfe, manages regional and high speed train tickets all over Spain and The Gurus harvested the price of tickets sold by Renfe. We are interested in trips between

Madrid and Barcelona and, for now, ask the question: are tickets more expensive one way or another? To answer this, we consider a sample of 10000 tickets, but restrict attention to AVE tickets sold at Promo rate. Our test statistic will again be the mean difference between the price (in euros) for a train ticket for Madrid–Barcelona (μ_1) and the price for Barcelona–Madrid (μ_2), i.e., $\mu_1 - \mu_2$. The null hypothesis is that there are no difference in price, so $\mathcal{H}_0 : \mu_1 - \mu_2 = 0$. We again use Welch test statistic for two samples.

```
# Library for manipulating data, including the pipe operator (%>%)
library(poorman)
# Load data
data(renfe, package = "hecstatmod")
head(renfe, n = 5)
```

```
## # A tibble: 5 x 7
##   price type    class    fare    dest          duration wday
##   <dbl> <fct>   <fct>   <fct>   <fct>          <dbl> <fct>
## 1 143.  AVE     Preferente Promo   Barcelona-Madrid    190 6
## 2 182.  AVE     Preferente Flexible Barcelona-Madrid    190 2
## 3  86.8 AVE     Preferente Promo   Barcelona-Madrid    165 7
## 4  86.8 AVE     Preferente Promo   Barcelona-Madrid    190 7
## 5  69.0 AVE-TGV Preferente Promo   Barcelona-Madrid    175 4
```

```
# Sub-sample with only Promo tickets
renfe_promo <- renfe %>% subset(fare == "Promo")
# two-sample t-test and mean difference
ttest <- t.test(price~dest, data = renfe_promo)
ttest #print result
```

```
##
## Welch Two Sample t-test
##
## data: price by dest
## t = -1, df = 8040, p-value = 0.2
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
## -1.100 0.209
## sample estimates:
## mean in group Barcelona-Madrid mean in group Madrid-Barcelona
##                                82.1                                82.6
```

Rather than use the asymptotic distribution, whose validity stems from the central limit theorem, we could consider another approximation under the less restrictive assumption that the data are exchangeable: under the null hypothesis, there is no difference between the two destinations and so the label for destination (a binary indicator) is arbitrary. The reasoning underlying permutation tests is as follows: to create a benchmark, we will consider observations with the same number in each group, but permuting the labels. We then compute the test statistic on each of these datasets. If there are only a handful in each group (fewer than 10), we could list all possible permutations of the data, but otherwise we can repeat this procedure many times, say 9999, to get a good approximation. This gives an approximate distribution from which we can extract the p -value by computing the rank of our statistic relative to the others.

```
# p-value (permutation test)
n <- nrow(renfe_promo)
B <- 1e4
ttest_stats <- numeric(B)
ttest_stats[1] <- ttest$statistic
set.seed(20200608) # set seed of pseudo-random number generator
for(i in 2:B){
  # Recalculate the test statistic, permuting the labels
  ttest_stats[i] <- t.test(price ~ dest[sample.int(n = n)],
                           data = renfe_promo)$statistic
}
# Graphics library
library(ggplot2)
# Plot the empirical permutation distribution
ggplot(data = data.frame(statistic = ttest_stats),
       aes(x=statistic)) +
  geom_histogram(bins = 30, aes(y=..density..), alpha = 0.2) +
  geom_density() +
  geom_vline(xintercept = ttest_stats[1]) +
  ylab("density") +
  stat_function(fun = dnorm, col = "blue")
```

The so-called bootstrap approximation to the p -value of the permutation test, 0.186, is the proportion of statistics that are more extreme than the one based on the original sample. It is nearly identical to that obtained from the Satterthwaite approximation, 0.182 (the Student- t distribution is numerically equivalent to a standard normal with that many degrees of freedom), as shown in Figure 1.5. Even if our sample is very large ($n = 8059$ observations), the difference is not statistically significant. With a bigger sample (the database has more

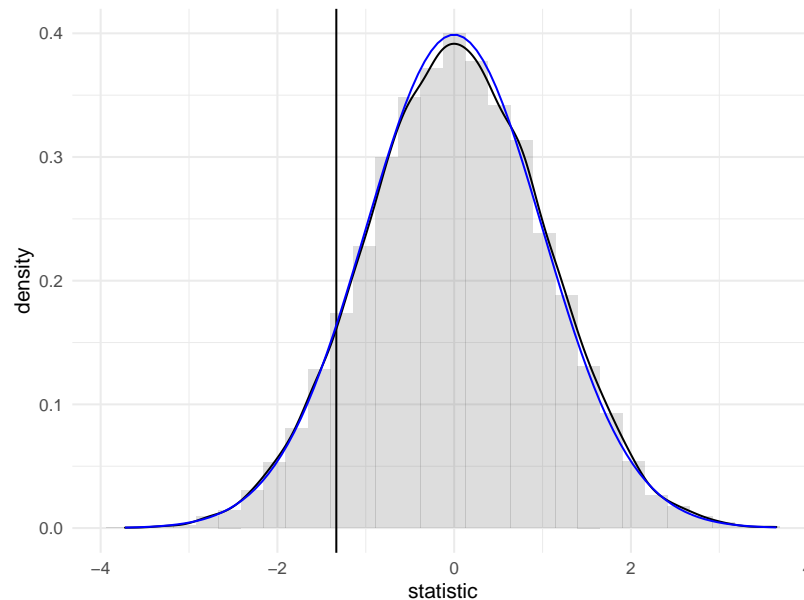


Figure 1.5: Permutation-based approximation to the null distribution of Welch two-sample t-test statistic (histogram and black curve) with standard normal approximation (blue curve) for the price of AVE tickets at promotional rate between Madrid and Barcelona. The value of the test statistic calculated using the original sample is represented by a vertical line.

than 2 million tickets), we could estimate more precisely the average difference, up to 1/100 of an euro: the price difference would eventually become statistically significant, but this says nothing about practical difference: 0.28 euros relative to an Promo ticket priced on average 82.56 euros is a negligible amount.

1.2 Exploratory Data Analysis

Before fitting a model, it is advisable to understand the structure of the data to avoid interpretation errors. Basic knowledge of graphs is required and we will spend some time addressing this. Further references include

- Chapter 3, *R for Data Science* by Garrett Golemund and Hadley Wickham
- Section 1.6 of OpenIntro *Introductory Statistics with Randomization and Simulation*
- *Fundamentals of Data Visualization* by Claus O. Wilke
- Chapter 1 of *Data Visualization: A practical introduction* by Kieran Healy

If exploratory data analysis is often neglected in statistics (perhaps because it has little to

no mathematical foundations), it is crucial. More than a rigorous approach, it is an art: Grolmund and Wickham talk of “state of mind”. The purpose of graphical exploratory data analysis is the extraction of useful information, often through a series of preliminary questions that are refined as the analysis progresses. Of particular interest are the relations and interactions between different variables and the distribution of the variables themselves. The major steps for undertaking an exploratory analysis are:

1. Formulate questions about the data
2. Look for answers using frequency table, descriptive statistics and graphics.
3. Refine the questions in light of the finding

In a report, you should highlight the most important features in a summary so that the reader can grasp your understanding and so that you guide him or her in the interpretation of the data.

1.2.1 Polish your work

Pay as much attention to figures and tables as to the main text. These should always include a legend that describes and summarizes the findings in the graph (so that the latter is standalone), name of variables (including units) on the axes, but also proper formatting so that the labels and numbers are readable (good printing quality, not too small). One picture is worth 1000 words, but make sure the graph tells a coherent story and that it is mentioned in the main text. Also ensure that only the necessary information is displayed: superfluous information (spurious digits, useless summary statistics) should not be presented.

1.2.2 Variable type

The data we will handle are stored in tables or frames. If the data frame is stocked in long format, each line corresponds to an observation and each column to a variable: the entries of the data base contain the (numeric) values.

The alternative is wide format, whereby the columns represent categorical variables and the entries are values of the response for a specific category (notably contingency tables). Figure 1.6 shows the difference between the two structures. Software typically require long formatted database for modelling purposes.

- a **variable** represents a characteristic of the population, for example the sex of an individual, the price of an item, etc.
- an **observation** is a set of measures (variables) collected under identical conditions for an individual or at a given time.

The choice of statistical model and test depends on the underlying type of the data collected. There are many choices: quantitative (discrete or continuous) if the variables are numeric,

wide				long		
id	x	y	z	id	key	val
1	a	c	e	1	x	a
2	b	d	f	2	x	b
				1	y	c
				2	y	d
				1	z	e
				2	z	f

Figure 1.6: Long versus wide-format for data tables (illustration by Garrick Aden-Buie).

or qualitative (binary, nominal, ordinal) if they can be described using an adjective; I prefer the term categorical, which is more evocative.

Most of the models we will deal with are so-called regression models, in which the mean of a quantitative variable is a function of other variables, termed explanatories. There are two types of numerical variables

- a discrete variable takes a countable number of values, prime examples being binary variables or count variables.
- a continuous variable can take (in theory) an infinite possible number of values, even when measurements are rounded or measured with a limited precision (time, width, mass). In many case, we could also consider discrete variables as continuous if they take enough values (e.g., money).

Categorical variables take only a finite of values. They are regrouped in two groups, nominal if there is no ordering between levels (sex, color, country of origin) or ordinal if they are ordered (Likert scale, salary scale) and this ordering should be reflected in graphs or tables. We will bundle every categorical variable using arbitrary encoding for the levels: for modelling, these variables taking K possible values (or levels) must be transformed into a set of $K - 1$ binary 0/1 variables, the omitted level corresponding to a baseline. Failing to declare categorical variables in your favorite software is a common mistake, especially when these are saved in the database using integers rather than strings.

1.2.3 Graphs

The main type of graph for representing categorical variables is bar plot (and modifications thereof). In a bar plot, the frequency of each category is represented in the y -axis as a function of the (ordered) levels on the x -axis. This representation is superior to the ignominious pie chart, a nuisance that ought to be banned (humans are very bad at comparing areas and a simple rotation changes the perception of the graph)!

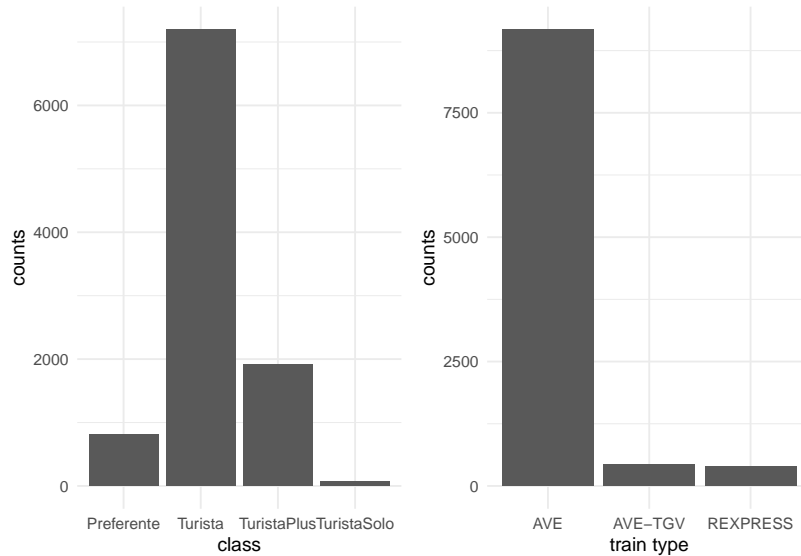


Figure 1.7: Bar plot of ticket class for Renfe tickets data

Continuous variables can take as many distinct values as there are observations, so we cannot simply count the number of occurrences by unique values. Instead, we bin them into distinct intervals so as to obtain an histogram. The number of class depends on the number of observations: as a rule of thumb, the number of bins should not exceed \sqrt{n} , where n is the sample size. We can then obtain the frequency in each class, or else normalize the histogram so that the area under the bands equals one: this yields a discrete approximation of the underlying density function. Varying the number of bins can help us detect patterns (rounding, asymmetry, multimodality).

Since we bin observations together, it is sometimes difficult to see where they fall. Adding rugs below or above the histogram will add observation about the range and values taken, where the heights of the bars in the histogram carry information about the (relative) frequency of the intervals.

If we have a lot of data, it sometimes help to focus only on selected summary statistics. A box-and-whiskers plot (or boxplot) represents five numbers

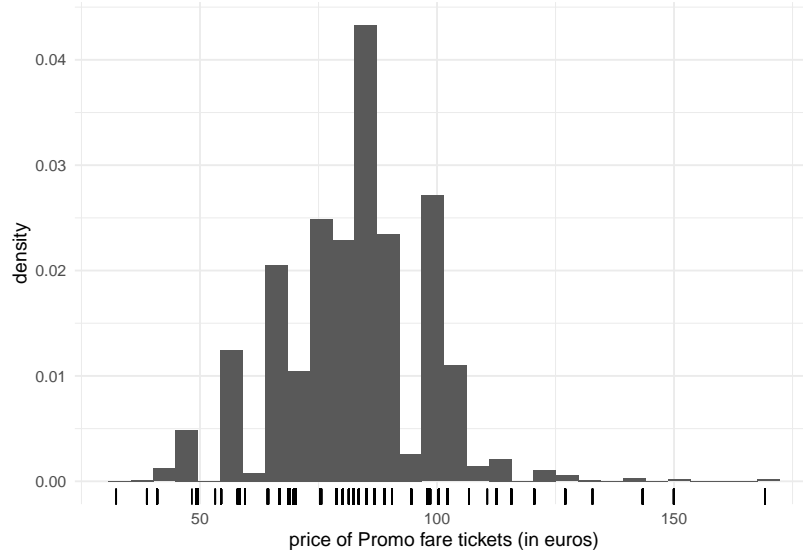


Figure 1.8: Histogram of Promo tickets for Renfe ticket data

- The box gives the quartiles q_1 , q_2 , q_3 of the distribution. The middle bar q_2 is thus the median, so 50% of the observations are smaller or larger than this number.
- The length of the whiskers is up to 1.5 times the interquartiles range $q_3 - q_1$ (the whiskers extend until the latest point in the interval, so the largest observation that is smaller than $q_3 + 1.5(q_3 - q_1)$, etc.)
- Observations beyond the whiskers are represented by dots or circles, sometimes termed outliers. However, beware of this terminology: the larger the sample size, the more values will fall outside the whiskers. This is a drawback of boxplots, which was conceived at a time where the size of data sets was much smaller than what is current standards.

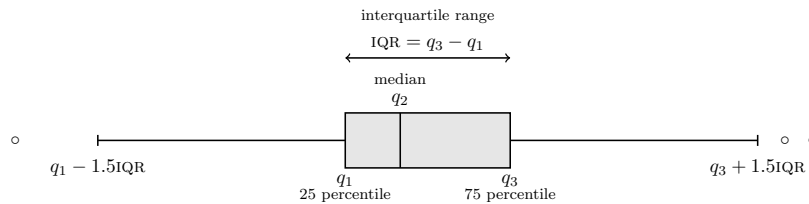


Figure 1.9: Box-and-whiskers plot

We can represent the distribution of a response variable as a function of a categorical variable by drawing a boxplot for each category and laying them side by side. A third

variable, categorical, can be added via a color palette, as shown in Figure 1.10.

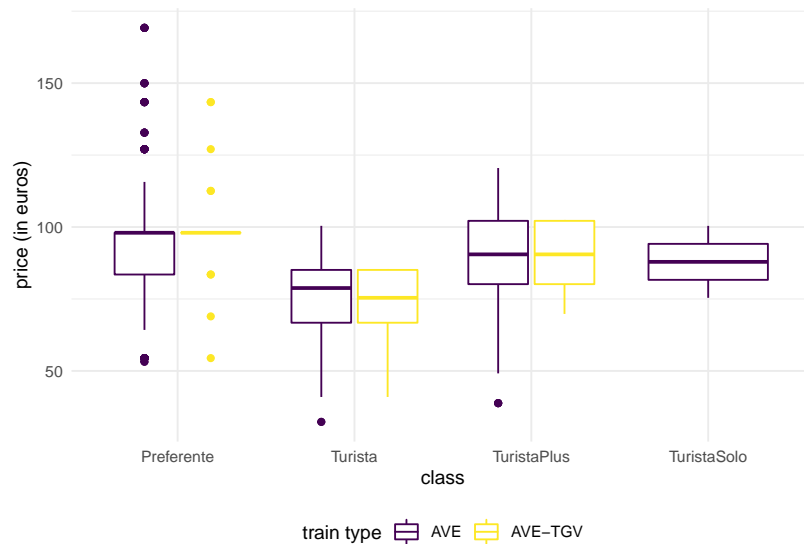


Figure 1.10: Box-and-whiskers plots for Promo fare tickets as a function of class and type for the Renfe tickets data.

Scatterplots are used to represent graphically the co-variation between two continuous variables: each tuple gives the coordinate of the point. If only a handful of large values are visible on the graph, a transformation may be useful: oftentimes, you will encounter graphs where the x - or y -axis is on the log-scale when the underlying variable is positive. If the number of data points is too large, it is hard to distinguish points because they are overlaid: adding transparency, or binning using a two-dimensional histogram with the frequency represented using color are potential solutions. The left panel of Figure 1.11 shows the 100 simulated observations, whereas the right-panel shows a larger sample of 10 000 points using hexagonal binning, an analog of the bivariate density.

Sometimes, continuous data have a particular structure, mostly when observations are collected over space or time. Time series are ordered and the response should be plotted on the y -axis as a function of time (on the x -axis). It is customary to draw segments between observations, but this display is sometimes misleading.

1.2.4 Exploratory data analysis

Rather than describe in details the exploratory analysis procedure, we proceed with an example that illustrates the process on the Renfe ticket dataset that was introduced previously.

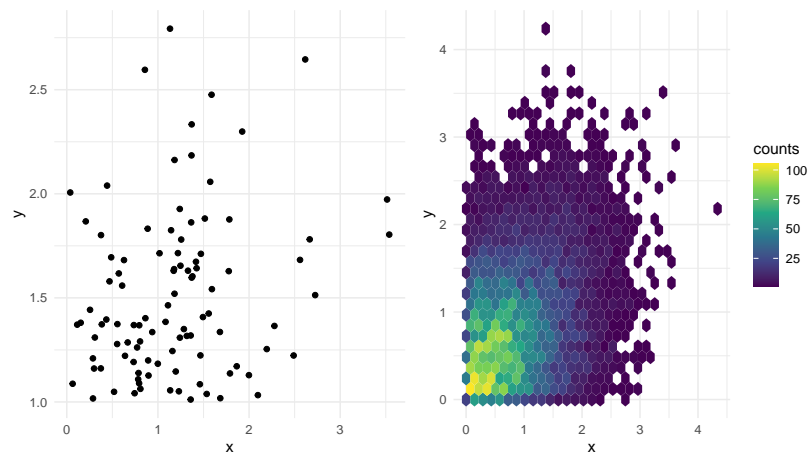


Figure 1.11: Scatterplot (left) and hexagonal heatmap of bidimensional bin counts (right) of simulated data.

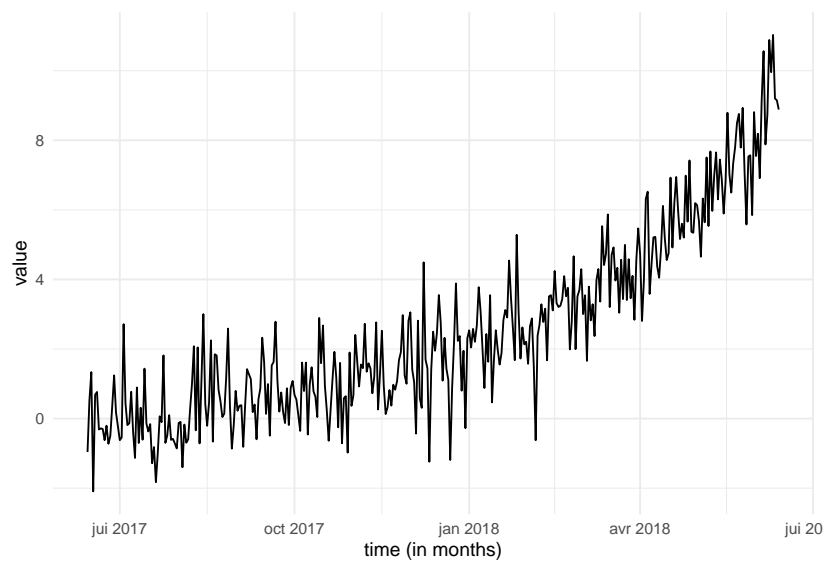


Figure 1.12: Graphical representation of a time series.

Example 1.3 (Exploratory data analysis of Renfe tickets). First, read the documentation accompanying the dataset! The data base `renfe` contains the following variables:

- `price` price of the ticket (in euros);
- `dest` binary variable indicating the journey, either Barcelona to Madrid (0) or Madrid to Barcelona (1);
- `fare` categorical variable indicating the ticket fare, one of `AdultoIda`, `Promo` or `Flexible`;
- `class` categorical variable giving the ticket class, either `Preferente`, `Turista`, `TuristaPlus` or `TuristaSolo`;
- `type` categorical variable indicating the type of train, either `Alta Velocidad Española (AVE)`, `Alta Velocidad Española jointly with TGV` (partnership between SNCF and Renfe for trains to/from Toulouse) `AVE-TGV` or regional train `REXPRESS`; only trains labelled `AVE` or `AVE-TGV` are high-speed trains.
- `duration` length of train journey (in minutes);
- `wday` categorical variable (integer) denoting the week day, ranging from Sunday (1) to Saturday (7).

There are no missing values and a quick view of the first row of the data frame (`head(renfe)`) shows that the data are stored in long format, meaning each line corresponds to a different ticket. We will begin our exploratory analysis with vague questions, for example

1. What are the factors determining the price and travel time?
2. Does travel time depend on the type of train?
3. What are the distinctive features of train types?
4. What are the main differences between fares?

Except for `price` and `duration`, all the other (explanatory) variables are categoricals. These need to be cast into factors (`factor`), especially integer-valued levels such as `wday`.

The database is clean and this preliminary preprocessing step has been done already. We can check the type of encoding using the command `str`, which also shows the data; the function `summary` is used to obtain descriptive statistics (min, max, mean, quartiles for continuous variables or else frequency for categorical variables); the function also returns the number of missing values (`NA`) of each column.

Data manipulation is often messy and **R** base syntax is particularly inelegant: data frames are list whose elements are accessed using `$`: for example `renfe$price`. A more legible and modular alternative is the pipe operator (`%>%`), with which one creates a logical chain of command (this function is not part of **R** base packages, but the libraries `tidyverse` and the minimal alternative `poorman` include it).

```
renfe %>% count(class)
```

```
##      class      n
## 1 Preferente  809
## 2   Turista 7197
## 3 TuristaPlus 1916
## 4 TuristaSolo   78
```

```
# `count` is a shortcut for the following syntax
renfe %>% group_by(type) %>% tally()
```

```
##      type      n
## 1     AVE 9174
## 2 AVE-TGV  429
## 3 REXPRESS 397
```

```
renfe %>% group_by(fare) %>% tally()
```

```
##      fare      n
## 1 AdultoIda  397
## 2 Flexible 1544
## 3   Promo 8059
```

By counting the number of train tickets in each category, we notice there are as many REXPRESS tickets as there are tickets sold at AdultoIda fare. Using a contingency table to get the number in respective sub-categories of each of those variables confirms that all tickets in the database for RegioExpress trains are sold with the AdultoIda fare and that there is only a single class, Turista. There are few such tickets, only 397 out of 10 000. This raises a new question: why are such trains so unpopular?

```
##      fare      type      n
## 1 AdultoIda REXPRESS  397
## 2 Flexible      AVE 1446
## 3 Flexible AVE-TGV   98
## 4   Promo      AVE 7728
## 5   Promo AVE-TGV  331
```

We have only scratched the surface, but one could also notice that there are only 17 duration values on tickets (`renfe %>% distinct(duration)` or `unique(renfe$duration)`). This leads us to think the duration on the ticket (in minutes) is the expected travel time. The

majority of those travel time (15 out of 17) are smaller than 3h15, but the other two exceed 9h! Looking at Google Maps, Madrid and Barcelona are 615km apart by car, 500km as the crow flies. this means some trains travel at about 200km/h, while others are closer to 70km/h. What are these slower trains? the variable `type` is the one most likely to encode this feature, and a quick look shows that the RegioExpress trains fall in the slow category (mystery solved!)

```
renfe %>%
  subset(duration > 200) %>%
  group_by(type, dest) %>%
  summarise("average duration" = mean(duration),
            "std. dev" = sd(duration),
            "average price" = mean(price),
            "std. dev" = sd(price))
```

##	type	dest	average duration	std. dev	average price	std. dev
## 1	REXPRESS	Barcelona-Madrid	544	0	43.2	0
## 2	REXPRESS	Madrid-Barcelona	562	0	43.2	0

The regular trains running between two cities take more than 9h, but one way (Madrid to Barcelona) is 18 minutes slower than in the other direction. More striking, we see that the price of the RegioExpress tickets is fixed: 43.25 euros regardless of direction. This is the most important finding so far, because these are not a sample for price: there is no variability! Graphics could have lead to the discovery (the boxplot of price as a function of train type would collapse to a single value).

We could have suspected that trains labeled AVE are faster: after all, it is the acronym of *Alta Velocidad Española*, literally Spanish high speed. What is the distinction between the two high speed train types. According to the SNCF website, AVE-TGV trains are partnership between Renfe and SNCF that operate between France and Spain.

```
renfe %>%
  subset(type %in% c("AVE", "AVE-TGV")) %>%
  group_by(type, dest) %>%
  summarise("average duration" = mean(duration),
            "std. dev" = sd(duration),
            "average price" = mean(price),
            "std. dev" = sd(price))
```

##	type	dest	average duration	std. dev	average price	std. dev
## 1	AVE	Barcelona-Madrid	171	15.9	87.4	19.8

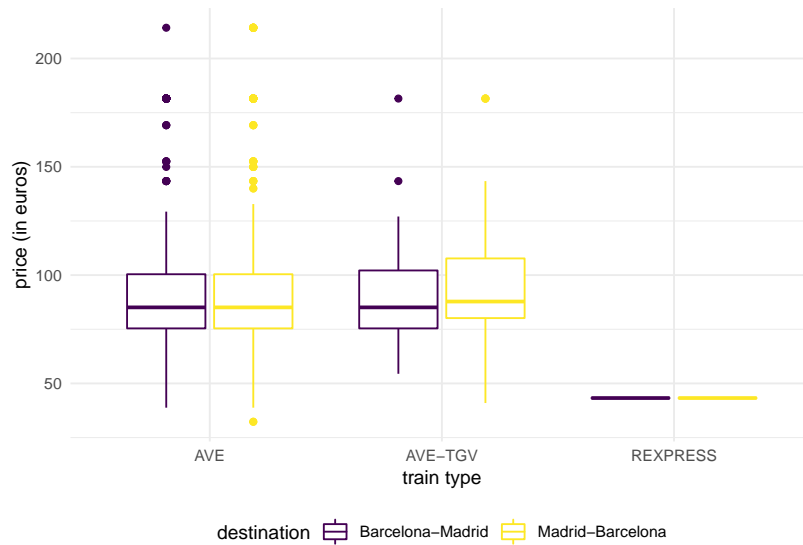


Figure 1.13: Boxplot of ticket price as a function of destination and train type.

## 2	AVE Madrid-Barcelona	170	16.6	88.2	20.8
## 3	AVE-TGV Barcelona-Madrid	175	0.0	87.0	16.8
## 4	AVE-TGV Madrid-Barcelona	179	0.0	90.6	20.2

The price of high speed trains are on average more than twice as expensive as regular trains. There is strong evidence of heterogeneity (standard deviation of 20 euros), which should raise scrutiny and suggests that high speed train tickets are dynamically priced. There is a single duration time for AVE-TGV tickets. We do not see meaningful differences in price depending on the type or the direction, but fares of ticket class availability may differ depending on whether the train is run in partnership with SNCF.

We have not looked at ticket fare and class, except for RegioExpress trains. Figure 1.15 shows large disparity in the variance of price according to fare: Promo fare takes many more distinct values than AdultoIda (duh) and Flexible fares. First class tickets (Preferente) are more expensive, but there are fewer observations falling in this group. Turista class is the least expensive for high-speed trains and the most popular. TuristaPlus offer an alternative to the latter with more comfort, whereas TuristaSolo gives access to individual seats.

Fare-wise Promo and PromoPlus give access to rebates that can go up to 70% and 65%, respectively. Promo tickets cannot be cancelled or exchanged, while both are possible with PromoPlus by paying a penalty amounting to 30-20% of the ticket price. Flexible fare ticket a sold at the same price as regular high-speed train tickets, but offer additional benefits

(and no rebates!)

```
renfe %>% subset(fare != "AdultoIda") %>%
ggplot(aes(y = price, x = class, col = fare)) +
  geom_boxplot() +
  labs(y = "price (in euros)",
       x = "class",
       color = "fare") +
  theme(legend.position = "bottom")
```

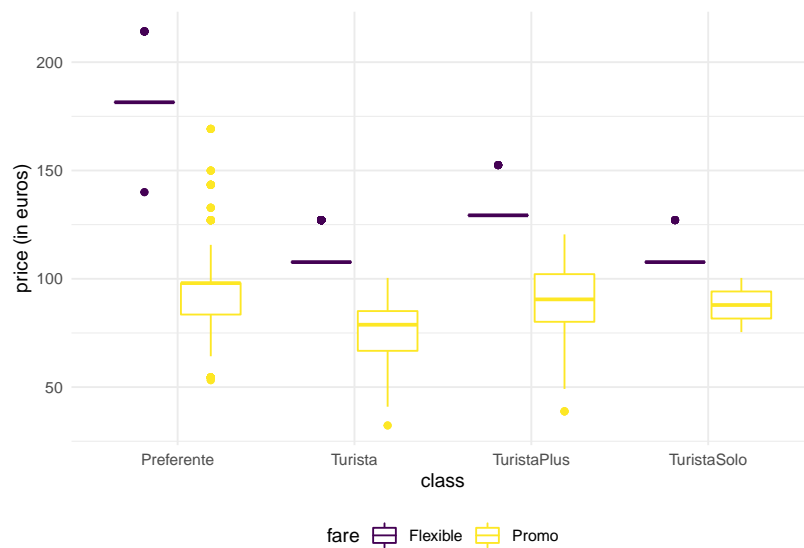


Figure 1.14: Boxplot of ticket price as a function of fare and class for high-speed Renfe trains.

```
ggplot(data = renfe, aes(x = price, y=..density.., fill = fare)) +
  geom_histogram(binwidth = 5) +
  labs(x = "price (in euros)", y = "density") +
  theme(legend.position = "bottom")
```

```
# Check the spread of Flexible tickets
renfe %>% subset(fare == "Flexible") %>% count(price, class)
```

```
##   price      class    n
```

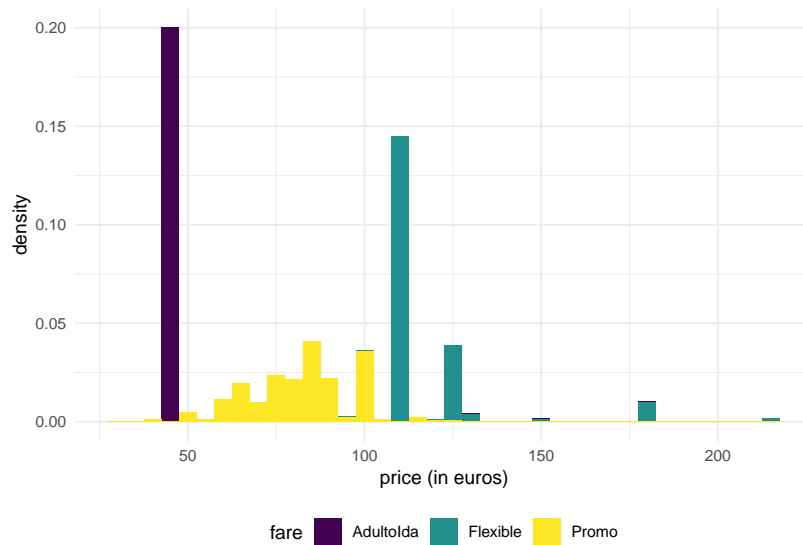


Figure 1.15: Histograms of ticket price as a function of fare for Renfe trains.

```
## 1  108      Turista 1050
## 2  108 TuristaSolo   67
## 3  127      Turista 285
## 4  127 TuristaSolo    9
## 5  129 TuristaPlus   31
## 6  140 Preferente    2
## 7  152 TuristaPlus   10
## 8  182 Preferente   78
## 9  214 Preferente   12
```

Note how Flexible tickets prices are spread: the boxplot is crushed and the interquartile range seems zero, even if some of the values are larger: this is indicative either constant price or of (too few) tickets in the category. We can find out which of these two possibilities is most likely by counting the number of Flexible fare tickets for the different types.

Neither duration, nor type or destination explains why some Flexible tickets are more or less expensive than the average. Promo tickets, on the other hand, are cheaper on average and Preferente more expensive.

We can summarize our findings:

- more than 91% of trains are high-speed trains.
- travel time depends on the type of train: high-speed train take at most 3h20.

- duration records expected travel time: there are only 17 unique values, 13 of which are for AVE trains.
- the price of RegioExpress train ticket is fixed (43.25€); all such tickets are sold with AdultoIda fare and there only one class (Turista). 57% of these trains go from Barcelona to Madrid and travel time is 9h22 from Barcelona to Madrid, 9h04 in the other direction.
- Turista is the cheapest and most popular class. Preferente class tickets are more expensive and are less often sold. TuristaPlus offers more comfort and TuristaSolo let you get individual seats.
- according to the Renfe website, Flexible fare tickets “come with additional offers and passengers can exchange or cancel their tickets if they miss their train”; as a counterpart, these tickets are more expensive and most tickets have a fixed fare (a handful are cheaper or more expensive, but this price difference is unexplained).
- the distribution of Promo fare high-speed trains ticket prices are more or less symmetric, but Flexible tickets seem left-truncated (the minimum price for these tickets in the sample is 107.7€).
- it appears that tickets sold by Renfe (Promo fare) are dynamically priced: the latter can be up to 70% cheaper than regular high-speed train tickets when purchased through the official agency or Renfe’s website. These tickets cannot be refunded or exchanged.
- there is no indication that prices depend on the direction of travel.

Chapter 2

Linear regression

A linear regression is a model for the conditional mean of a response variable Y as a function of p explanatory variables (also termed regressors or covariates),

$$E(Y \mid \mathbf{X}) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p. \quad (2.1)$$

The mean of Y is conditional on the values of the observed covariate \mathbf{X} ; this amounts to treating them as non-random, known in advance.

In practice, any model is an approximation of reality. An error term is included to take into account the fact that no exact linear relationship links \mathbf{X} and Y (otherwise this wouldn't be a statistical problem), or that measurements of Y are subject to error. The random error term ε will be the source of information for our inference, as it will quantify the goodness of fit of the model.

We can rewrite the linear model in terms of the error for a random sample of size n : denote by Y_i the value of the response for observation i , and X_{ij} the value of the j th explanatory variable of observation i . The model is

$$Y_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip} + \varepsilon_i, \quad i = 1, \dots, n, \quad (2.2)$$

where ε_i is the additive error term specific to observation i . While we may avoid making distributional assumption about ε_i , we nevertheless fix its expectation to zero to encode the fact we do not believe the model is systematically off, so $E(\varepsilon_i \mid \mathbf{X}_i) = 0$ ($i = 1, \dots, n$).

One important remark is that the model is linear in the coefficients $\beta \in \mathbb{R}_{p+1}$, not in the explanatory variables! the latter are arbitrary and could be (nonlinear) functions of other explanatory variables, for example $X = \ln(\text{years})$, $X = \text{horsepower}^2$ or $X = I_{\text{man}} \cdot I_{\text{full}}$. The mean of the response is specified as a **linear combination of explanatory variables**. This is at the core of the flexibility of the linear regression, which is used mainly for the following purposes:

1. Evaluate the effects of covariates \mathbf{X} on the mean response of Y .
2. Quantify the influence of the explanatories \mathbf{X} on the response and test for their significance.
3. Predict the response for new sets of explanatories \mathbf{X} .

2.1 Introduction

Linear regression is the most famous and the most widely used statistical model around. The name may appear reductive, but many tests statistics (t -tests, ANOVA, Wilcoxon, Kruskal–Wallis) can be formulated using a linear regression, while models as diverse as trees, principal components and deep neural networks are just linear regression model in disguise. What changes under the hood between one fancy model to the next are the optimization method (e.g., ordinary least squares, constrained optimization or stochastic gradient descent) and the choice of variables entering the model (spline basis for non-parametric regression, indicator variable selected via a greedy search for trees, activation functions for neural networks).

This chapter explores the basics of linear regression, parameter interpretation and testing for coefficients and sub-models. Analysis of variance will be presented as special case of linear regression.

To make concepts and theoretical notions more concrete, we will use data from a study performed in a college in the United States. The goal of the administration who collected these information was to investigate potential gender inequality in the salary of faculty members. The data contains the following variables:

- salary: nine-month salary of professors during the 2008–2009 academic year (in thousands USD).
- rank: academic rank of the professor (assistant, associate or full).
- field: categorical variable for the field of expertise of the professor, one of applied or theoretical.
- sex: binary indicator for sex, either man or woman.
- service: number of years of service in the college.
- years: number of years since PhD.

Before drafting a model, a quick look at the data is in due order. If salary increases with year, there is more heterogeneity in the salary of higher ranked professors: logically, assistant professors are either promoted or kicked out after at most 6 years according to the data. The limited number of years prevents large variability for their salaries.

Salary increases over years of service, but its variability also increases with rank. Note the much smaller number of women in the sample: this will impact our power to detect differences between sex. A contingency table of sex and academic rank can be useful to see

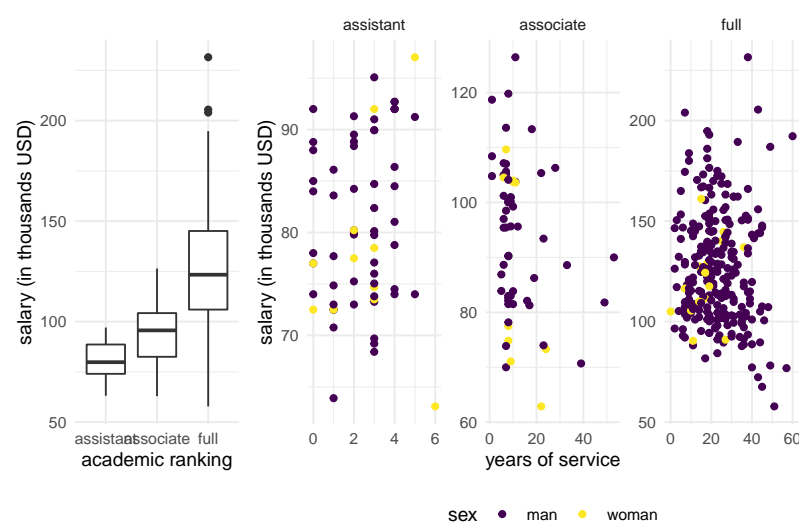


Figure 2.1: Exploratory data analysis of college data: salaries of professors as a function of the number of years of service and the academic ranking

if the proportion of women is the same in each rank: women represent 16% of assistant professors and 16% of associate profs, but only 7% of full professors and these are better paid on average.

Contingency table of the number of prof in the college by sex and academic rank.

assistant

associate

full

man

56

54

248

woman

11

10

18

The simple linear regression model only includes a single explanatory variable and defines a straight line linking two variables X and Y by means of an equation of the form $y = a + bx$; Figure 2.2 shows the line passing through the scatterplot for years of service.

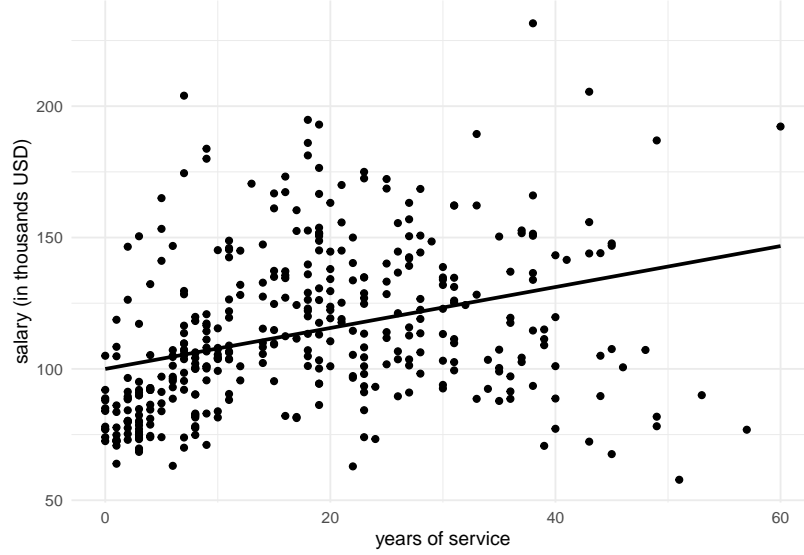


Figure 2.2: Simple linear regression model for the salary of professors as a function of the number of years of service; the line is the solution of the least squares problem.

2.2 Ordinary least squares

The ordinary least square estimators $\hat{\beta} = (\hat{\beta}_0, \dots, \hat{\beta}_p)$ are the values that simultaneously minimize the Euclidean distance between the random observations Y_i and the **fitted values**

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_{i1} + \dots + \hat{\beta}_p X_{ip}, \quad i = 1, \dots, n.$$

In other words, the least square estimators are the solution of the convex optimization problem

$$\hat{\beta} = \min_{\beta \in \mathbb{R}^{p+1}} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2 = \min_{\beta} \|Y - X\beta\|^2$$

This system of equations has an explicit solution which is better expressed using matrix notation: this amounts to expressing equation (2.2) with one observation per line.

Consider the matrices

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}, \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}, \quad \mathbf{X} = \begin{pmatrix} X_{11} & X_{12} & \cdots & X_{1p} \\ X_{21} & X_{22} & \cdots & X_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ X_{n1} & X_{n2} & \cdots & X_{np} \end{pmatrix}, \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_1 \\ \beta_2 \\ \vdots \\ \beta_p \end{pmatrix}$$

The model in compact form is

$$\mathbf{Y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}.$$

The ordinary least squares estimator solves the unconstrained optimization problem

$$\hat{\boldsymbol{\beta}} = \min_{\boldsymbol{\beta} \in \mathbb{R}^{p+1}} (\mathbf{y} - \mathbf{X}\boldsymbol{\beta})^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}).$$

and a proof is provided in the Appendix. If the $n \times (p+1)$ matrix \mathbf{X} is full-rank, we obtain a unique solution to the optimization problem,

$$\hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{Y}. \quad (2.3)$$

What does the solution to the least squares problem represent in two dimensions? The estimator is the one minimizing the sum of squared residuals: the i th **ordinary residual** $e_i = y_i - \hat{y}_i$ is the *vertical* distance between a point y_i and the fitted value \hat{y}_i on the line; the blue segments on Figure 2.3 represent the individual vectors of residuals.

Remark (Geometry of least squares). If we consider the n observations as a (column) vector, the term $\mathbf{X}\hat{\boldsymbol{\beta}}$ is the projection of the response vector \mathbf{y} on the linear span generated by the columns of \mathbf{X} , $\mathcal{S}_{\mathbf{X}}$. The ordinary residuals are thus orthogonal to $\mathcal{S}_{\mathbf{X}}$ and to the fitted values, meaning $\mathbf{e}^\top \hat{\mathbf{y}} = 0$. A direct consequence of this fact is that the linear correlation between \mathbf{e} and $\hat{\mathbf{y}}$ is zero; we will use this property to build graphical diagnostics.

Remark (Complexity of ordinary least squares). This is an aside: in machine learning, you will often encounter linear models fitted using a (stochastic) gradient descent algorithm. Unless your sample size n or the number of covariates p is significant (think at the Google scale), an approximate should not be preferred to the exact solution! From a numerical perspective, obtaining the least square estimates requires inverting a $(p+1) \times (p+1)$ matrix $\mathbf{X}^\top \mathbf{X}$, which is the most costly operation. In general, direct inversion should be avoided because it is not the most numerically stable way of obtaining the solution. **R** uses the QR decomposition, which has a complexity of $O(np^2)$. Another more stable alternative, which has the same complexity but is a bit more costly is use of a singular value decomposition.

Any good software will calculate ordinary least square estimates for you. Keep in mind that there is an explicit and unique solution provided your design matrix \mathbf{X} doesn't contain collinear columns. If you have more than one explanatory variable, the fitted values lie on a hyperplan (which is hard to represent graphically). Mastering the language and technical term (fitted values, ordinary residuals, etc.) is necessary for the continuation.

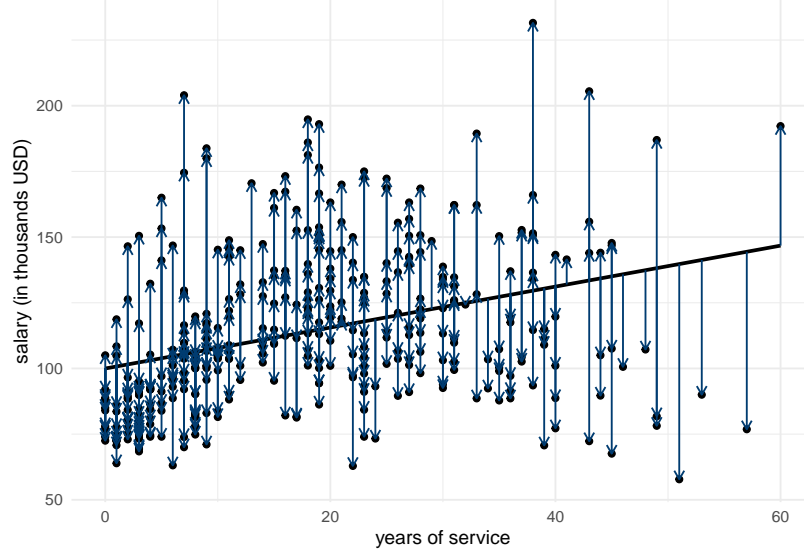


Figure 2.3: Illustration of ordinary residuals added to the regression line (blue vectors).

2.3 Interpretation of the model parameters

What do the β parameters of the linear model represent? In the simple case presented in Figure 2.2, the equation of the line is $\hat{Y} = \hat{\beta}_0 + \hat{\beta}_1 X_1$, β_0 is the intercept (the mean value of Y when $X_1 = 0$) and β_1 is the slope, i.e., the average increase of Y when X_1 increases by one unit.

In some cases, the intercept is **meaningless** because the value $X_1 = 0$ is impossible (e.g., X_1 represents the height of a human). In the same vein, there may be no observation in a neighborhood of $X_1 = 0$, even if this value is plausible, in which case the intercept is an **extrapolation**.

If the columns of \mathbf{X} are arbitrary, it is customary to include an intercept: this amounts to including $\mathbf{1}_n$ as column of the design matrix \mathbf{X} . Because the residuals are orthogonal to the columns of \mathbf{X} , their mean is zero, since $n^{-1} \mathbf{1}_n^\top \mathbf{e} = \bar{e} = 0$. In general, we could also obtain mean zero residuals by including a set of vectors in \mathbf{X} that span $\mathbf{1}_n$.

In our example, the equation of the fitted line of Figure 2.2 is

$$\widehat{\text{salary}} = 99.975 + 0.78 \text{service}.$$

The average salary of a new professor would then be 99974.653 dollars, whereas the average annual increase for each additional year of service is 779.569 dollars.

If the response variable Y should be continuous (for the least square criterion to be meaningful), we place no such restriction on the explanatories. The case of dummies in particular is common: these variables are encoded using binary indicators (0/1). Consider for example the sex of the professors in the study:

$$\text{sex} = \begin{cases} 0, & \text{for men,} \\ 1, & \text{for women.} \end{cases}$$

The equation of the simple linear model that includes the binary variable `sex` is $\text{salary} = \beta_0 + \beta_1 \text{sex} + \varepsilon$. Let μ_0 denote the average salary of men and μ_1 that of women. The intercept β_0 can be interpreted as usual: it is the average salary when $\text{sex} = 0$, meaning that $\beta_0 = \mu_0$. We can write the equation for the conditional expectation for each sex,

$$E(\text{salary} \mid \text{sex}) = \begin{cases} \beta_0, & \text{sex} = 0 \text{ (men),} \\ \beta_0 + \beta_1 & \text{sex} = 1 \text{ (women).} \end{cases}$$

A linear model that only contains a binary variable X as regressor amounts to specifying a different mean for each of two groups: the average of women is $E(\text{salary} \mid \text{sex} = 1) = \beta_0 + \beta_1 = \mu_1$ and $\beta_1 = \mu_1 - \mu_0$ represents the difference between the average salary of men and women. The least square estimator $\hat{\beta}_0$ is the sample mean of men and $\hat{\beta}_1$ is the difference of the sample mean of women and men. The parametrization of the linear model with β_0 and β_1 is in terms of **contrasts** and is particularly useful if we want to test for mean difference between the groups, as this amounts to testing $\mathcal{H}_0 : \beta_1 = 0$. If we wanted our model to directly output the sample means, we would need to replace the design matrix $\mathbf{X} = [\mathbf{1}_n, \text{sex}]$ by $[\mathbf{1}_n - \text{sex}, \text{sex}]$. The fitted model would be the same because they span the same 2D subspace, but this is not recommended because software treat cases without intercept differently and it can lead to unexpected behaviour (more on this latter).

If we fit the model with `sex` only to the `college` data, we find that the average salary of men is $\hat{\beta}_0 = 1.151 \times 10^5$ USD and the mean difference estimate of the salary between women and men is $\hat{\beta}_1 = 14088.009$ dollars. Since the estimate is negative, this means women are paid less. Bear in mind that the model is not adequate for determining if there are gender inequalities in the salary distribution: 2.2 shows that the number of years of service and the academic rank strongly impact wages, yet the distribution of men and women is not the same within each rank.

Even if the linear model defines a line, the latter is only meaningful when evaluated at 0 or 1; Figure 2.4 shows it in addition to sample observations (jittered) and a density estimate for each sex. The colored dot represents the mean, showing that the line does indeed pass through the mean of each group.

A binary indicator is a categorical variable with two levels, so we could extend our reasoning and fit a model with a categorical explanatory variable with k levels. To do this, we add

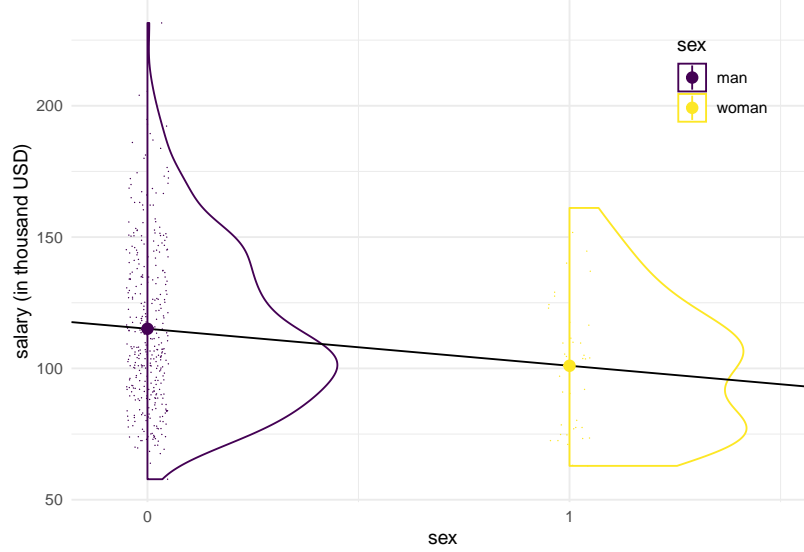


Figure 2.4: Simple linear model for the `college` data using the binary variable `sex` as regressor: even if the equation defines a line, only its values in 0/1 are realistic.

$k - 1$ indicator variables plus the intercept: if we want to model a different mean for each of the k groups, it is logical to only include k parameters in the mean model. We will choose, as we did with `sex`, a **reference category** or **baseline** whose average will be encoded by the intercept β_0 . The other parameters $\beta_1, \dots, \beta_{k-1}$ are contrasts relative to the baseline. The `college` data includes the ordinal variable `rank`, which has three levels (assistant, associate and full). We thus need two binary variables, $X_1 = \mathbb{I}(\text{rank} = \text{associate})$ and $X_2 = \mathbb{I}(\text{rank} = \text{full})$; the i th element of the vector X_1 is one for an associate professor and zero otherwise. The linear model is

$$\text{salary} = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \varepsilon,$$

and the conditional expectation of salary

$$E(\text{salary} \mid \text{rank}) = \begin{cases} \beta_0, & \text{rank} = \text{assistant}, \\ \beta_0 + \beta_1 & \text{rank} = \text{associate}, \\ \beta_0 + \beta_2 & \text{rank} = \text{full}, \end{cases}$$

Thus β_1 (respectively β_2) are the difference between the average salary of associate (respectively full) professors and assistant professors. The choice of the baseline category is arbitrary and all choices yield the same model: only the interpretation changes from one parametrization to the next. For an ordinal variable, it is recommended to choose the smallest or the largest category to ease comparisons.

The models we have fitted so far are not adequate because they ignore variables that are necessary to correctly explain variations in salaries: Figure 2.1 show for example that rank is critical for explaining the salary variations in the college. We should thus fit a model that include those simultaneously to investigate the gender gap (which consists of differences that are unexplained by other factors). Before doing this, we come back to the interpretation of the parameters in the multiple linear regression setting.

Consider the model $Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon$. The intercept β_0 represents the mean value of Y when *all* of the covariates are set to zero,

$$\beta_0 = E(Y \mid X_1 = 0, X_2 = 0, \dots, X_p = 0).$$

For categorical variables, this yields the baseline, whereas we fix the continuous variables to zero: again, this may be nonsensical depending on the study. The coefficient β_j ($j \geq 1$) can be interpreted as the mean increase of the response Y when X_j increases by one unit, all other things being equal (*ceteris paribus*); e.g.,

$$\begin{aligned} \beta_1 &= E(Y \mid X_1 = x_1 + 1, X_2 = x_2, \dots, X_p = x_p) \\ &\quad - E(Y \mid X_1 = x_1, X_2 = x_2, \dots, X_p = x_p) \\ &= \{\beta_0 + \beta_1(x_1 + 1) + \beta_2 x_2 + \dots + \beta_p x_p\} \\ &\quad - \{\beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p\} \end{aligned}$$

It is not always possible to fix the value of an explanatory if multiple columns of X contains functions/transformations of it. For example, if we included a polynomial of order k for some variable X ,

$$Y = \beta_0 + \beta_1 X + \beta_2 X^2 + \dots + \beta_k X^k + \varepsilon.$$

If we include a term of order k , X^k , we must *always* include the lower order terms $1, X, \dots, X^{k-1}$ to make sure the resulting model is interpretable (otherwise, it amounts to a particular class of polynomials with some zero coefficients). Interpreting nonlinear effects (even polynomials, for which $k \leq 3$ in practice), is complicated because the effect of an increase of one unit of X *depends of the value of the latter*.

Example 2.1 (Auto data). We consider a linear regression model for the fuel autonomy of cars as a function of the power of their motor (measured in horsepower) from the auto dataset. The postulated model,

$$\text{mpg}_i = \beta_0 + \beta_1 \text{horsepower}_i + \beta_2 \text{horsepower}_i^2 + \varepsilon_i,$$

includes a quadratic term. Figure 2.5 shows the scatterplot with the fitted regression line, above which the line for the simple linear regression for horsepower is added.

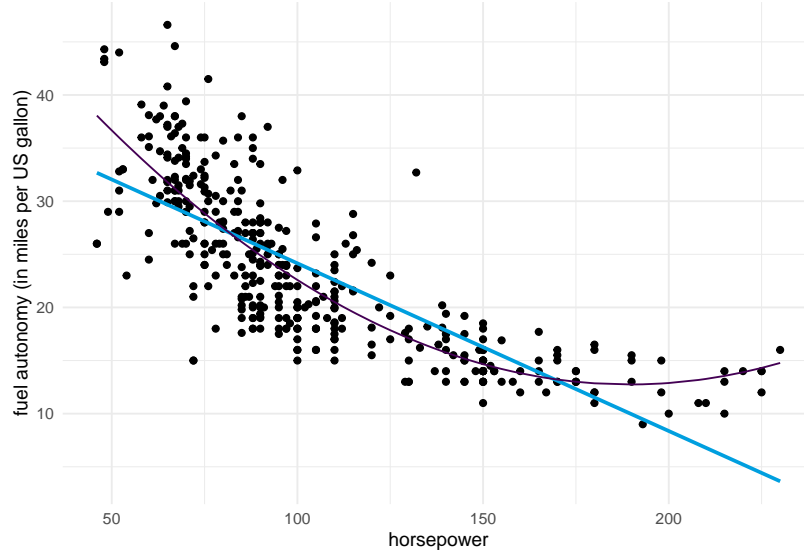


Figure 2.5: Linear regression models for the fuel autonomy of cars as a function of motor power

It appears graphically that the quadratic model fits better than the simple linear alternative: we will assess this hypothesis formally later. For the degree two polynomial, Figure 2.5 show that fuel autonomy decreases rapidly when power increases between 50 to 100, then more slow until 189.35 hp. After that, the model postulates that autonomy increases again as evidenced by the scatterplot, but beware of extrapolating (weird things can happen beyond the range of the data, as exemplified by Hassett's cubic model for the number of daily cases of Covid19 in the USA).

The representation in Figure 2.5 may seem counter-intuitive given that we fit a linear model, but it is a 2D projection of 3D coordinates for the equation $\beta_0 + \beta_1 x - y + \beta_2 z = 0$, where $x = \text{horsepower}$, $z = \text{horsepower}^2$ and $y = \text{mpg}$. Physics and common sense force $z = x^2$, and so the fitted values lie on a curve in a 2D subspace of the fitted plan, as shown in grey in the 3D Figure 2.6.

Remark (There are better alternatives to polynomials for modelling nonlinear effects). Generally speaking, one uses flexible basis vectors (splines) rather than polynomials for smoothing when the relation between the response Y and an explanatory variable X is nonlinear; these models involve many covariates and it is customary to add a penalty term to control for overfitting and wiggleness. A better (physical) understanding of the system, or a theoretical model can also guide the choice of functions to use.

The coefficient β_j in Eq. (2.1) represents the *marginal* contribution of X_j when all the other

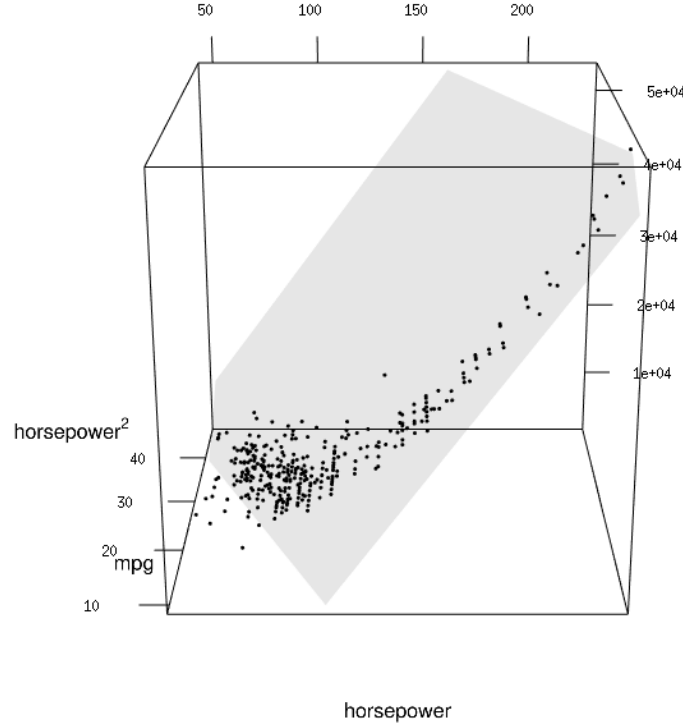


Figure 2.6: 3D graphical representation of the linear regression model for the `auto` data.

covariates are included in the model and which is not explained by them. This can be represented graphically by projecting Y and X_j in the orthogonal complement of \mathbf{X}_{-j} (the matrix containing all but the j th column X_j). The **added-variable plot** is a graphical tool showing this projection: the residuals from the linear regression of Y onto $\mathcal{S}(\mathbf{X}_{-j})$ are mapped to the y -axis, whereas the residuals from the linear regression of X_j as a function of \mathbf{X}_{-j} are shown on the x -axis. The regression line passes through $(0, 0)$ and its slope is $\hat{\beta}_j$. This graphical diagnostic is useful for detecting collinearity and the impact of outliers.

Example 2.2 (Wage inequality in an American college). We consider a multiple regression model for the `college` data that includes `sex`, `academic rank`, `field of study` and the number of years of service as regressors.

If we multiply `salary` by a thousand to get the resulting estimates in US dollars, the postulated model is

$$\begin{aligned} \text{salary} \times 1000 = & \beta_0 + \beta_1 \text{sex}_{\text{woman}} + \beta_2 \text{field}_{\text{theoretical}} \\ & + \beta_3 \text{rank}_{\text{associate}} + \beta_4 \text{rank}_{\text{full}} + \beta_5 \text{service} + \varepsilon. \end{aligned}$$

Estimated coefficients of the linear model for the `college` (in USD, rounded to the nearest dollar).

$$\hat{\beta}_0$$

$$\hat{\beta}_1$$

$$\hat{\beta}_2$$

$$\hat{\beta}_3$$

$$\hat{\beta}_4$$

$$\hat{\beta}_5$$

86596

-4771

-13473

14560

49160

-89

The interpretation of the coefficients is as follows:

- The estimated intercept is $\hat{\beta}_0 = 86596$ dollars; it corresponds to the mean salary of men assistant professors who just started the job and works in an applied domain.
- everything else being equal (same field, academic rank, and number of years of service), the estimated salary difference between a woman and is estimated at $\hat{\beta}_1 = -4771$ dollars.
- *ceteris paribus*, the salary difference between a professor working in a theoretical field and one working in an applied field is β_2 dollars: our estimate of this difference is -13473 dollars, meaning applied pays more than theoretical.
- *ceteris paribus*, the estimated mean salary difference between associate and assistant professors is $\hat{\beta}_3 = 14560$ dollars.
- *ceteris paribus*, the estimated mean salary difference between full and assistant professors is $\hat{\beta}_4 = 49160$ dollars.
- au sein d'un même échelon, chaque année supplémentaire de service mène à une augmentation de salary annuelle moyenne de $\hat{\beta}_5 = -89$ dollars.

All other factors taken into account, women get paid less than men. It remains to see if this difference is statistically significant. Perhaps more surprising, the model estimates that salary decreases with every additional year of service: this seems counterintuitive when looking at Figure 2.2, which showed a steady increase of salary over the years. However, this

graphical representation is misleading because Figure 2.1 showed that academic ranking mattered the most. Once all the other factors are accounted for, years of service serves to explain the salary of full professors who have been working unusual amounts of time and who gain less than the average full professor, as shown by the added-variable plot of Figure 2.7.

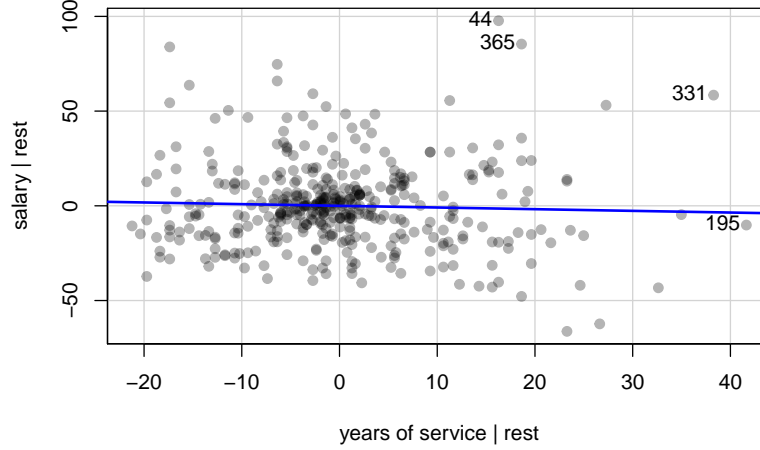


Figure 2.7: Added-variable plot for years of service in the linear regression model of the college data.

Details about implementation of linear models using **R** are provided in the Appendix.

2.4 Tests for parameters of the linear model

We need to make further assumptions in order to carry out inference and build testing procedures for the mean model parameters of the linear model. In order to get a tractable distribution for test statistics, it is customary to assume that the disturbances ε are independent normal random variables with mean zero and common variance σ^2 . It follows that Y_1, \dots, Y_n are conditionally *independent* random variables with

$$E(Y_i | \mathbf{X}_i) = \beta_0 + \sum_{j=1}^p X_{ij}\beta_j, \quad \text{Va}(Y_i | \mathbf{X}_i) = \sigma^2, \quad i = 1, \dots, n.$$

Under this hypothesis, the least square estimators for the mean parameters coincide with the maximum likelihood estimators. The advantage of imposing this (more stringent than

necessary) assumption is that we can use our toolbox for testing. The theory underlying likelihood tests is presented in the chapter on likelihood-based inference. Assuming normal errors leads to exact distributions for the test statistics (which also coincide with the asymptotic ones in large samples).

Of particular interest are tests of restrictions for components of β . The large sample properties of the maximum likelihood estimator imply that

$$\hat{\beta} \sim \text{No}_{p+1} \left\{ \beta, \sigma^2 (\mathbf{X}^\top \mathbf{X})^{-1} \right\}$$

for sufficiently large sample size. One can thus easily estimate the standard errors from the matrix upon replacing σ^2 by an estimator, typically the unbiased estimator of the variance.

In an inferential setting, it's often important to test whether the effect of an explanatory variable is significant: if X_j is binary or continuous, the test for $\mathcal{H}_0 : \beta_j = 0$ corresponds to a null marginal effect for X_j . The null model is a linear regression in which we remove the $(j + 1)$ st column of \mathbf{X} , so both models are nested. The Wald test statistic is reported by most software $W = \hat{\beta}_j / \text{se}(\hat{\beta}_j)$ and the null distribution is Student- t with $n - p - 1$ degrees of freedom, which explains the terminology (t values). In addition to coefficient estimates, it is possible to obtain confidence intervals for β_j , which are the usual $\hat{\beta}_j \pm t_{n-p-1, \alpha/2} \text{se}(\hat{\beta}_j)$, with $t_{n-p-1, \alpha/2}$ denoting the $1 - \alpha/2$ quantile of the St_{n-p-1} distribution.

For categorical variables with more than two levels, testing if $\beta_j = 0$ is typically not of interest because the contrast represent the different between the category X_j and the baseline: these two categories could have a small difference, but the categorical variable as a whole may still be a useful predictor given the other explanatories. The hypothesis of zero contrast is awkward because it implies a null model in which selected categories are merged.

2.4.1 F -tests for comparison of nested linear models

Consider the *full* linear model which contains p predictors,

$$\mathbb{M}_1 : Y = \beta_0 + \beta_1 X_1 + \dots + \beta_g X_g + \beta_{k+1} X_{k+1} + \dots + \beta_p X_p + \varepsilon.$$

Suppose without loss of generality that we want to test $\mathcal{H}_0 : \beta_{k+1} = \beta_{k+2} = \dots = \beta_p = 0$ (one could permute columns of the design matrix to achieve this configuration). The global hypothesis specifies that $(p - k)$ of the β parameters are zero. The *restricted model* corresponding to the null hypothesis contains only the covariates for which $\beta_j \neq 0$,

$$\mathbb{M}_0 : Y = \beta_0 + \beta_1 X_1 + \dots + \beta_k X_k + \varepsilon.$$

Let $\text{SS}_e(\mathbb{M}_1)$ be the residuals sum of squares for model \mathbb{M}_1 ,

$$\text{SS}_e(\mathbb{M}_1) = \sum_{i=1}^n (Y_i - \hat{Y}_i^{\mathbb{M}_1})^2,$$

where $\hat{Y}_i^{\mathbb{M}_1}$ is the i th fitted value from \mathbb{M}_1 . Similarly define $SS_e(\mathbb{M}_0)$ for the residuals sum of square of \mathbb{M}_0 . Clearly, $SS_e(\mathbb{M}_0) \geq SS_e(\mathbb{M}_1)$ (why?)

The F -test statistic is

$$F = \frac{\{SS_e(\mathbb{M}_0) - SS_e(\mathbb{M}_1)\}/(p - k)}{SS_e(\mathbb{M}_1)/(n - p - 1)}$$

Under \mathcal{H}_0 , the F statistic follows a Fisher distribution with $(p - k)$ and $(n - p - 1)$ degrees of freedom, $F(p - k, n - p - 1) - p - k$ is the number of restrictions, $n - p - 1$ is sample size minus the number of β 's in \mathbb{M}_1 .

It turns out that both F and t -statistics are equivalent for testing a single coefficient β_j : the F -statistic is the square of the Wald statistic and they lead to the same inference — the p -value for the test are identical. While it may get reported in tables, the test for $\beta_0 = 0$ is not of interest; we keep the intercept merely to centre the residuals.

For normal linear regression, the likelihood ratio test for comparing models \mathbb{M}_1 and \mathbb{M}_0 is a function of the sum of squared residuals: the usual formula simplifies to $R = n \ln\{SS_e(\mathbb{M}_0)/SS_e(\mathbb{M}_1)\}$. This is reminiscent of the F -statistic formula and the two are in fact intimately related modulo null distribution and scaling. The t -tests and F -tests presented above could thus both be viewed as particular cases of likelihood-based tests.

Consider the college data example and the associated linear model with rank, sex, years of service and field as covariates.

Table of linear regression coefficients with associated standard errors, Wald tests and p -values based on Student- t distribution

term

estimate

std. error

Wald stat.

p-value

(Intercept)

86.596

2.96

29.25

< 0.001

sex [woman]

-4.771
3.878
-1.23
0.22
field [theoretical]
-13.473
2.315
-5.82
< 0.001
rank [associate]
14.56
4.098
3.55
< 0.001
rank [full]
49.16
3.834
12.82
< 0.001
service
-0.089
0.112
-0.8
0.43

Table 2.4.1 shows the estimated coefficients (in thousands of dollars). The coefficients are the least squares estimates $\hat{\beta}$, the standard errors are the square root of the diagonal elements of $S^2(\mathbf{X}^\top \mathbf{X})^{-1}$. The Wald (or t -) statistic is simply $W = \hat{\beta}/\text{se}(\hat{\beta})$ for $\mathcal{H}_0 : \beta_j = 0$: given two of the three estimates, we could easily recover the third using the formula for the test. The p -values are for the two-sided alternative test with $\mathcal{H}_a : \beta_j \neq 0$.

The interpretation is usual: p -values that are less than our prescribed level α do not contribute significantly given the other variables already in the model. Neither years of service nor sex are statistically different from zero given all the other variables. The test for β_{sex} is comparing the model with all covariates (including service), and vice-versa. Note that the conclusion changes depending on the model: both coefficients would be statistically significant had we removed rank from the set of covariates, because they are correlated. The gender imbalance among ranks explains most of the gap between sex, whereas year of service is largely redundant once we account for the jumps due to change of academic rank.

Type 3 sum of square decomposition table: sum of square decomposition comparing nested models with and without covariates, F -statistic and associated p -value.

variable

sum of square

df

F stat.

p-value

(Intercept)

439059.2

1

855.71

< 0.001

sex

776.7

1

1.51

0.22

field

17372.5

1

33.86

< 0.001

```

rank
102883.1
2
100.26
< 0.001
service
324.5
1
0.63
0.43
Residuals
200620.4
391

```

Table 2.4.1 gives the F -statistics values and the associated p -values. The sum of squares represent the difference $SS_e(\mathbb{M}_0) - SS_e(\mathbb{M}_1)$ for various null models \mathbb{M}_0 , except the last line for residuals which reports $SS_e(\mathbb{M}_1)$. You can verify that (up to rounding) these p -values are identical to those of the Wald test in the output when $df=1$. The only categorical variable here with more than one level is rank, and it is strongly significant: removing it from the model leads to a sharp decrease in fit.

We could have also easily computed the likelihood ratio test to compare the models: for example, the log likelihood for the full model is -1799.027 and that of the model without rank is -1881.202, so the likelihood ratio statistic would be 164.349 and this is strongly significant when compared to the χ^2_2 distribution (both likelihood ratio test and F -test give a p -value of 2.2×10^{-16}).

2.5 Coefficient of determination

When we specify a model, the error term ε accounts for the fact no perfect linear relationship characterizes the data (if it did, we wouldn't need statistic to begin with). Once we have fitted a model, we estimate the variance σ^2 ; one may then wonder which share of the total variance in the sample is explained by the model.

The total sum of squares, defined as the sum of squared residuals from the intercept-only model, serves as comparison — the simplest model we could come up with would involving

every observation by the sample mean of the response and so this gives (up to scale) the variance of the response, $SS_c = \sum_{i=1}^n (y_i - \bar{y})^2$. We can then compare the variance of the original data with that of the residuals from the model with covariate matrix \mathbf{X} , defined as $SS_e = \sum_{i=1}^n e_i^2$ with $e_i = y_i - \hat{\beta}_0 - \sum_{j=1}^p \hat{\beta}_j X_{ij}$. We define the coefficient of determination, or squared multiple correlation coefficient of the model, R^2 , as

$$R^2 = 1 - \frac{SS_e}{SS_c} = \frac{\sum_{i=1}^n (y_i - \bar{y})^2 - \sum_{i=1}^n e_i^2}{\sum_{i=1}^n (y_i - \bar{y})^2}.$$

The coefficient of determination can be interpreted as the square of Pearson's linear correlation between the response y and the fitted values \hat{y} ; see the Appendix for a derivation of this fact.

It's important to note that R^2 is not a goodness-of-fit criterion: some phenomena are inherently noisy and even a good model will fail to account for much of the response's variability. Moreover, one can inflate the value of R^2 by including more explanatory variables: the coefficient is non-decreasing in the dimension of \mathbf{X} , so a model with $p+1$ covariate will necessarily have a higher R^2 value than only p of the explanatories. For model comparisons, it is better to employ information criteria, or else rely on the predictive performance if this is the purpose of the regression. Lastly, a model with a high R^2 may imply high correlation, but the relation may be spurious: linear regression does not yield causal models!

2.6 Predictions

When we compute least square estimates, we obtain fitted values \hat{y} as $\mathbf{X}\hat{\beta}$, where \mathbf{X} denotes the $n \times (p+1)$ matrix of original observations. Recalling that $E(Y_i | \mathbf{X}_i) = \beta_0 + \sum_{j=1}^p \beta_j X_{ij}$, we can obtain an estimate of the mean surface for any value of $\mathbf{x}_{n+1} \in \mathbb{R}^p$ by replacing the unknown coefficients β by our estimates $\hat{\beta}$ – this actually yields the best linear unbiased predictor of the mean.

If we want to predict the value of a new observation, say Y_{n+1} , with explanatories $\mathbf{x}_{n+1} = (1, x_1, \dots, x_p)$, the prediction of the value will also be $\hat{y}_{n+1} = \mathbf{x}_{n+1}\hat{\beta}$ because

$$E(Y_{n+1} | \mathbf{x}_{n+1}) = \mathbf{x}_{n+1}\beta + E(\varepsilon_{n+1} | \mathbf{x}_{n+1}) = \mathbf{x}_{n+1}\beta.$$

However, individual observations vary more than averages (which are themselves based on multiple observations). Intuitively, this is due to the added uncertainty of the error term ε_{n+1} appearing in the model equation: the variability of new predictions is the sum of uncertainty due to the estimators (based on random data) and the intrinsic variance of the observations assuming the new observation is independent of those used to estimate the

coefficients,

$$\begin{aligned} \text{Va}(Y_{n+1} \mid \mathbf{x}_n) &= \text{Va}(\mathbf{x}_{n+1}\hat{\boldsymbol{\beta}} + \varepsilon_{n+1} \mid \mathbf{x}_n) \\ &= \text{Va}(\mathbf{x}_{n+1}\hat{\boldsymbol{\beta}} \mid \mathbf{x}_n) + \text{Va}(\varepsilon_{n+1} \mid \mathbf{x}_n) \\ &= \sigma^2 \mathbf{x}_{n+1}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}^\top + \sigma^2, \end{aligned}$$

where S^2 is the unbiased estimator of the variance σ^2 . The distribution of Y_{n+1} is by assumption normal, but since we do not know the variance, we base the prediction interval on the Student distribution, viz.

$$\frac{\mathbf{x}_{n+1}\hat{\boldsymbol{\beta}} - Y_{n+1}}{\sqrt{S^2\{1 + \mathbf{x}_{n+1}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}^\top\}}} \sim \text{St}_{n-p-1}.$$

and obtain $1 - \alpha$ prediction interval for Y_{n+1} by inverting the test statistic,

$$\mathbf{x}_{n+1}\hat{\boldsymbol{\beta}} \pm t_{n-p-1}(\alpha/2) \sqrt{S^2\{1 + \mathbf{x}_{n+1}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}^\top\}}.$$

Similar calculations yield the formula for confidence intervals for the mean,

$$\mathbf{x}_{n+1}\hat{\boldsymbol{\beta}} \pm t_{n-p-1}(\alpha/2) \sqrt{S^2 \mathbf{x}_{n+1}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_{n+1}^\top}.$$

The two differ only because of the additional variance of individual observations.

Figure 2.8 shows pointwise uncertainty bands for a simple linear regression of the `intention` data as a function of `fixation`, illustrating the limitations of the linear model in this example: the model is not accounting for the fact that our response arises from a bounded discrete distribution with integer values ranging from 2 to 14. The middle line gives the prediction of individuals as we vary their fixation time. Looking at the formula for the confidence, it is clear that the bands are not linear (we consider the square root of a function that involves the predictors), but it is not obvious that the uncertainty increases as you move away from the average of the predictors. This is more easily seen by replicating the potential curves that could have happened with different data: I generated new potential slopes from the asymptotic normal distribution of $\hat{\boldsymbol{\beta}}$ estimators to show the hyperbolic shape is not surprising: we are basically tilting curves from the average fixation/intention, and they have higher potential from deviating far from the range of observations.

2.7 Interactions

In the interpretation of the linear model, the effect of an explanatory variable is assumed to be the same, regardless of the other explanatory variables (*ceteris paribus*). To isolate the effect of β_j , we indeed fix the value of the other explanatories and increase by one the

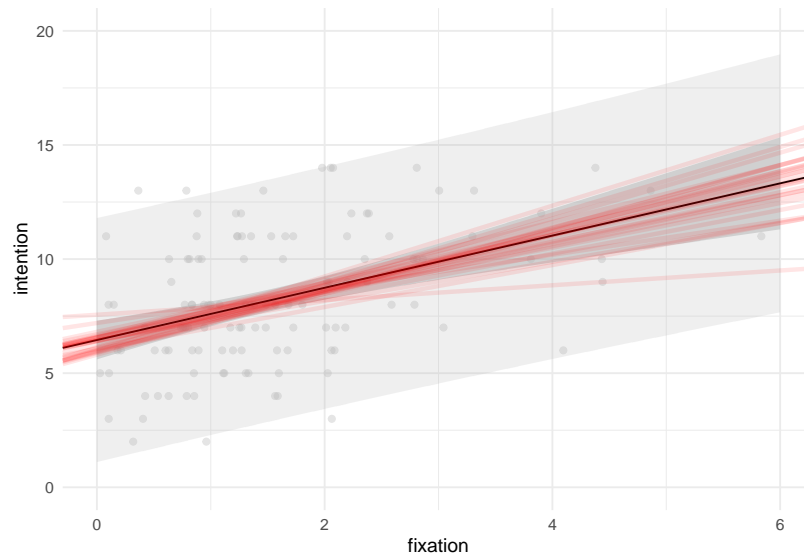


Figure 2.8: Prediction for the simple linear regression of buying intention as a function of fixation time. The plot shows predictions along with pointwise 95% confidence intervals of the mean and the individual predictions.

variable X_j (whenever this makes sense) to obtain the slope coefficient. However, the effect of a covariate may sometimes depend on another explanatory.

A good example of interaction is provided by the insurance dataset. An exploratory data analysis suggested that premiums depended on age, smoker status and body mass index, but through obesity status. It can be best represented graphically by looking at body mass index.

From there, we could create an indicator variable $obese = I(bmi \geq 30)$ and add an interaction term between smoker/obese (categorical) and age (continuous) in our mean model. We take non-smoker as baseline category. To make interpretation more meaningful, we rescale age so that $age = 0$ corresponds to 18 years old.

Table of regression coefficients for the insurance data with interaction terms between age, obesity and smoker status.

term

estimate

std. error

Wald stat.



Figure 2.9: Graph of insurance charges against body mass index. The figure clearly shows the interaction: the premium vary for smokers depending on whether or not they are obese, but we see no such behaviour for non-smokers. There is a clear linear increase of charges with age, but it is not clear whether the annual increase is the same for the three groups.

p-value

(Intercept)

2668.7

362.7

7.36

< 0.001

age

265.9

15

17.7

< 0.001

obesity [obese]

115.3

510.3

0.23

0.82

smoker [yes]

13526.2

803.1

16.84

< 0.001

age * obesity [obese]

1.4

20.1

0.07

0.94

age * smoker [yes]

-5.3

33.5

-0.16

0.87

obesity [obese] * smoker [yes]

19308.7

1110.5

17.39

< 0.001

age * obesity [obese] * smoker [yes]

19.1

44.8

0.43

0.67

The linear regression model has eight parameters, which could be mapped to four intercepts and four different slopes for age; however, the model is parametrized in terms of contrasts, which facilitates testing restrictions.

$$\begin{aligned} \text{charges} = & \beta_0 + \beta_1 \text{age} + \beta_2 \text{obese} + \beta_3 \text{smoker} + \beta_4 \text{age} \cdot \text{obese} \\ & + \beta_5 \text{age} \cdot \text{smoker} + \beta_6 \text{obese} \cdot \text{smoker} + \beta_7 \text{age} \cdot \text{obese} \cdot \text{smoker} + \varepsilon \end{aligned}$$

Because of the three-way interaction, it is not possible to recover individual parameters by changing the value of the corresponding covariate and keeping everything else constant: changing the smoker status likely impacts multiple regressors simultaneously. To retrieve the interpretation of the different coefficients, we will need to change one parameter at the time, write the mean equation and then isolate the coefficients. Throughout, obese is a dummy variable equal to one if the person has a body mass index greater than 30 and likewise smoker if the person is a smoker.

$$\text{charges} = \begin{cases} \beta_0 + \beta_1 \text{age} & (g_1) \text{ non-obese, non-smoker} \\ (\beta_0 + \beta_2) + (\beta_1 + \beta_4) \text{age}, & (g_2) \text{ obese, non-smoker} \\ (\beta_0 + \beta_3) + (\beta_1 + \beta_5) \text{age}, & (g_3) \text{ non-obese, smoker} \\ (\beta_0 + \beta_2 + \beta_3 + \beta_6) + (\beta_1 + \beta_4 + \beta_5 + \beta_7) \text{age}, & (g_4) \text{ obese, smoker} \end{cases}$$

- The intercept β_0 is the average at 18 years old of non-smokers who are not obese.
- The slope β_1 is the average annual increase in charges for non-smokers who are not obese.
- The parameter β_2 is a contrast, the difference between the average charges of 18 years old non-smokers who are obese and those who are not.
- The parameter β_3 is a contrast, the difference between the average premium for non-obese 18 years old who smoke and those who don't.
- The parameter β_4 is a contrast, the difference in average annual increase for non-smokers between obese and non-obese adults.
- The parameter β_5 is a contrast, the difference in average annual increase for non-obese between smoker and non-smoker adults.

The other two coefficients, β_6 and β_7 represent differences of average between groups $g_1 + g_4 - g_2 - g_3$ for both intercepts and slopes.

The only F -test that we should consider in the analysis of variance table (containing the Type III sum of squares decomposition) is the test for the three-way interaction, which corresponds to $\mathcal{H}_0 : \beta_7 = 0$. The p -value against the two-sided alternative is 0.6704, suggesting no difference in slope. The reason why we do not consider the other tests is that they correspond to irrelevant hypotheses. For example, the test for the two-way interaction term $\mathcal{H}_0 : \beta_4 = 0$ associated to `age · obese` would correspond to merging the intercepts of group `g1` and `g2`. Changing the baseline category would imply that a different difference in intercept is forced to be zero.

Sometimes, however, specific hypothesis could be of interest because of the problem setting. We could perform bespoke test to check here that $\mathcal{H}_0 : \beta_2 = \beta_4 = 0$, which consists in merging the two non-smoker categories, or even $\mathcal{H}_0 : \beta_2 = \beta_4 = \beta_5 = \beta_7 = 0$, which amounts to merging non-smokers and the imposing a common slope for `age`. Such tests are not directly available in the output, but we can implement them manually by fitting two models and then plug-in in the values of residual sum of squares of both model in the formula of the F statistic. Using the F null distribution, we get a p -value of 0.965. This suggests both non-smoker groups are indistinguishable and likewise that there is no evidence that slopes for `age` are not equal.

Most interactions include functions of categorical variables together with either other categorical variables (different intercepts / conditional means per subgroups) or else other. Keep in mind that the validity of our tests above depend on the model being correctly specified: there is however evidence of difference in heterogeneity between groups, with unexplained non-smoker records. Plotting the residuals from the model that includes four different intercepts and slopes for `age` for each combination of smoker/obesity status misses other features that one would capture in diagnostic plots. In particular, except for some notable outliers, there is evidence that the premiums of smokers also increase with body mass index; as evidenced by Figure 2.10. If we include `bmi` as additional covariate in the model in addition of `obese`, the interpretation of changes in obesity will depend on the value of `bmi` and vice-versa.

2.7.1 Two-way ANOVA

“We’re the same as regression, but we’ve established a separate persona.” —
Chelsea Parlett-Pelleriti

Two-way analysis of variance is a linear regression model with two categorical variables, and possibly an interaction term between the two. We consider a study on the effect of the type of delay and source of delay on the evaluation of service and on time lost waiting.

Hui, M. K., Thakor, M. V. et Gill, R. (1998). *The Effects of Delay Type and Service Stage on Consumers’ Reaction to Waiting*. *Journal of Consumer Research* **24**, 469-479.

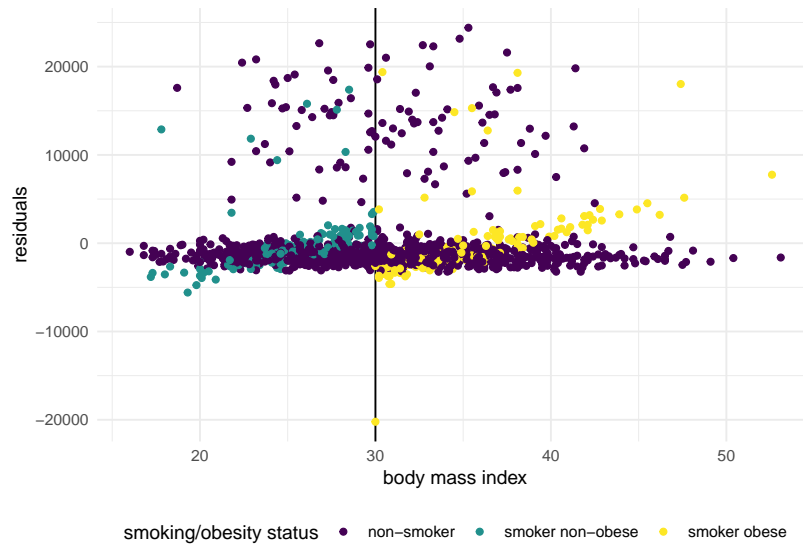


Figure 2.10: Residuals from the interaction model for charges with age, smoker and obesity. There is a notably outlier for a male smoker whose bmi is exactly 30 and other points above. There is indication of a linear trend for both smoker sub-groups as their body mass index increase, which was not apparent previously because of the effect of age.

In a university, 120 participants tried a new course registration system. They were actually using a fake registration system where factors of interest could be manipulated, the advancement stage in the registration process and the type of delay. The two levels for stage of advancement were close to end and far from end. In stage far from end, a message indicating a delay appeared immediately at the beginning of the registration process. For the stage close to end, the delay message appeared after the participant has entered personal information and course choices.

There were three types of delays: procedural, correctional, and unknown. The delay message indicated that the server was currently busy (procedural delay), that there were problems with the terminal and that the system needed to reestablish the connection (correctional delay) or else simply stated “please wait” (unknown delay).

At the end of the registration process, the participants were asked to estimate the delay time (in minutes) incurred during the registration process. They were also asked to provide an evaluation of the service using two measurement scales. The simulated data corresponding to this study are found in the `delay` database, which contains the following variables

- `time`: delay time (in minutes) according to the participant.
- `eval`: evaluation of service (standardized score).

- stage: stage of advancement, a factor with levels `close to end` and `far from end`
- delay: type of delay, a factor with levels `procedural`, `correctional` and `unknown`.

In the experimental design, the 120 participants were randomly assigned to one of these six conditions, but nine of the 120 participants were removed because they were not able to specify the type of delay that occurred. The dataset is **unbalanced**, meaning that the number of observations in each cell is unequal. If there had been the same number of observations in each subgroup for factors A , B , the test for the effect of factor A given that B is already in the model would be the same as the marginal effect of A only.

We will evaluate the effect of the factors `delay` and `stage` on the estimated waiting time.

If there is an interaction between explanatory variables, but the latter is not included in the models, the tests are misleading: one assumption of the linear model being that all of the relevant covariates have been included and their effect properly accounted for. In this example, omitting the interaction leads to effects being averaged and cancelling each other. Fitting an ANOVA model with only additive (main effects, without interaction) for the evaluation of service suggests that the factors `stage` and `delay` are not significant: the p -values for the main effects, reported in Table 2.7.1, are 0.409 and 0.137. Therefore, there seemingly is no significant difference in delay time between the two stages and the three delay types.

Analysis of variance table (Type 3 sum of square decomposition table) for the model for evaluation of service without interaction. While both factor appear not to be significant given the other, interaction plots show this is due to cancellation effects. The conclusions of the test are invalid because model assumptions are violated.

variable

sum of square

df

F stat.

p-value

(Intercept)

0

1

0.11

0.74

stage

0.3

1

0.69

0.41

delay

1.6

2

2.02

0.14

Residuals

41.9

105

However, looking at the data and the mean of each sub-class paints a different portrait, as evidenced by the interaction plot of Figure 2.11. The score for service increases for correctional delays when the error occurs close to the end, but the effect is opposite for other delay types. The fact that our p -values for the effects were not significant merely indicate that the effects cancelled out each other.

There is a significant interaction between the variables `delay` and `stage` on evaluation. Therefore, the effect of `delay` on the variable `eval` depends on the level of the variable `stage` and vice-versa. Since there is an interaction term, tests for the main effects of `stage` or `delay` are not of interest.

Analysis of variance table (Type 3 sum of square decomposition table). The only test of interest is that for the interaction, which is highly significant.

variable

sum of square

df

F stat.

p-value

(Intercept)

3.8

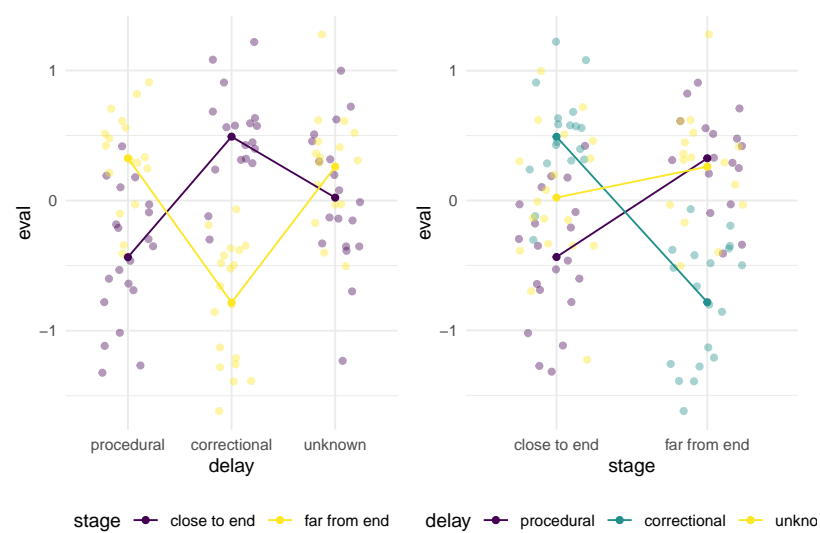


Figure 2.11: Interactions plots for the delay data, with jittered observations: if there was no interaction, the curves should be parallel. These plots seem to indicate an interaction.

1
18.29
< 0.001
stage
5.3
1
25.67
< 0.001
delay
8.1
2
19.64
< 0.001
stage:delay

20.5

2

49.63

< 0.001

Residuals

21.3

103

If a difference is detected at this stage, one could then proceed to compute all pairwise differences within a given level of the other categorical variable. This is equivalent to running multiple t -tests, but the linear model assumes equal variance and would thus pool all observations to estimate the latter. For example, we could test if the difference between procedural and correctional delays when the error occurs far from the end of the procedure is significant. The mean differences of a factor for a given level of another factor are called **simple effects**. If the interaction is significant, there must be at least one pairwise difference which is also significant. However, repeated testing can inflate the Type I error, as illustrated in this comic strip: if we conduct 20 independent tests and the null is true in all cases (so there are no difference between any of the sub-groups), we would still expect to reject (by mistake) $1/20$ on average for tests performed at level α . There exists method to adjust for multiple testing.

While the normality assumption is less crucial for the reliability of the tests, it requires sufficient sample size in each subgroup for these (say 20), so that the central limit theorem kicks in. The tests and models typically use the pooled variance estimators: in small samples, it makes better use of the data and has higher precision than estimators of the individual variance in each subgroups. If the variance were unequal, comparisons would require extended models covered later in the term and pairwise comparisons could be based on Welch test statistic.

2.8 Collinearity

The linearity assumption can be interpreted broadly to mean that all relevant covariates have been included and that their effect is correctly specified in the equation of the mean. Adding superfluous covariates to a model has limited impact: if the (partial) correlation between a column vector \mathbf{X}_k and the response variable Y is zero, then $\beta_k = 0$ and the estimated coefficient $\hat{\beta}_k \approx 0$ because the least square estimators are unbiased. If we include many useless variables, say k , the lack of parsimony can however make interpretation more

difficult. The price to pay for including the k additional covariates is an increase in the variance of the estimators $\hat{\beta}$.

It is nevertheless preferable to include more variables than to forget key predictors: if we omit an important predictor, their effect may be picked up by other regressors (termed **confounders**) in the model with are correlated with the omitted variable. The interpretation of the other effects can be severely affected by confounders. For example, the simple linear model (or two-sample t -test) for salary as a function of sex for the college data is invalid because sex is a confounder for rank. Since there are more men than women full professor, the mean salary difference between men and women is higher than it truly is. One way to account for this is to include control variables (such as rank), whose effect we need not be interested in, but that are necessary for the model to be adequate. We could also have used stratification, i.e., tested for wage discrimination within each academic rank. This is the reason why sociodemographic variables (sex, age, education level, etc.) are collected as part of studies.

A linear model is not a causal model: all it does is capture the linear correlation between an explanatory variable and the response. When there are more than one explanatory, the effect of X_j given what has not already been explained by X_{-j} . Thus, if we fail to reject $\mathcal{H}_0 : \beta_j = 0$ in favor of the alternative $\mathcal{H}_1 : \beta_j \neq 0$, we can only say that there is no significant *linear* association between X_j and Y once the effect of other variables included in the model has been accounted for. There are thus two scenarios: either the response is uncorrelated with X_j (uninteresting case, but easy to pick up by plotting both or computing linear correlation), or else there is a strong correlation between X_j and both the response Y as well as (some) of the other explanatory variables X_1, \dots, X_p . This problem is termed (multi)collinearity.

One potential harm of collinearity is a decrease in the precision of parameter estimators. With collinear explanatories, many linear combinations of the covariates represent the response nearly as well. Due to the (near) lack of identifiability, the estimated coefficients become numerically unstable and this causes an increase of the standard errors of the parameters. The predicted or fitted values are unaffected. Generally, collinearity leads to high estimated standard errors and the regression coefficients can change drastically when new observations are included in the model, or when we include or remove explanatories. The individual β coefficients may not be statistically significant, but the global F -test will indicate that some covariates are relevant for explaining the response. This however would also be the case if there are predictors with strong signal, so neither is likely to be useful to detect issues.

The added-variable plot shows the relation between the response Y and an explanatory X_j after accounting for other variables: the slope $\hat{\beta}_j$ of the simple linear regression is the same of the full model. A similar idea can be used to see how much of X_j is already explained by the other variables. For a given explanatory variable X_j , we define its **variance inflation**

factor as $VIF(j) = (1 - R^2(j))^{-1}$, where $R^2(j)$ is the coefficient of determination of the model obtained by regressing X_j on all the other explanatory variables, i.e.,

$$X_j = \beta_0^* + \beta_1^* X_1 + \cdots + \beta_{j-1}^* X_{j-1} + \beta_{j+1}^* X_{j+1} + \cdots + \beta_p^* X_p + \varepsilon^*$$

By definition, $R^2(j)$ represents the proportion of the variance of X_j that is explained by all the other predictor variables. Large variance inflation factors are indicative of problems (typically covariates with $VIF > 10$ require scrutiny, and values in the hundreds or more indicate serious problems).

Added-variable plots can also serve as diagnostics, by means of comparison of the partial residuals with a scatterplot of the pair (Y, X_j) ; if the latter shows very strong linear relation, but the slope is nearly zero in the added-variable plot, this hints that collinearity is an issue.

What can one do about collinearity? If the goal of the study is to develop a predictive model and we're not interested in the parameters themselves, then we don't need to do anything. Collinearity is not a problem for the overall model: it's only a problem for the individual effects of the variables. Their joint effect is still present in the model, regardless of how the individual effects are combined.

If we are interested in individual parameter estimates, for example, to see how (and to what extent) the predictor variables explain the behaviour of Y , then things get more complicated. Collinearity only affects the variables that are strongly correlated with one another, so we only care if it affects one or more of the variables of interest. There sadly is no good solution to the problem. One could

- try to obtain more data, so as to reduce the effects of collinearity appearing in specific samples or that are due to small sample size.
- create a composite score by somehow combining the variables showing collinearity.
- remove one or more of the collinear variables. You need to be careful when doing this not to end up with a misspecified model.
- use penalized regression. If $\mathbf{X}^\top \mathbf{X}$ is (nearly) not invertible, this may restore the uniqueness of the solution. Penalties introduce bias, but can reduce the variance of the estimators β . Popular choices include ridge regression (with an l_2 penalty), lasso (l_1 penalty), but these require adjustment in order to get valid inference.

Whatever the method, it's important to understand that it can be very difficult (and sometimes impossible) to isolate the individual effect of a predictor variable strongly correlated with other predictors.

Example 2.3 (Collinearity in the `college` data). We consider the `college` data analysis and include all the covariates in the database, including `years`, the number of years since PhD. One can suspect that, unless a professor started his or her career elsewhere before moving to the college, they will have nearly the same years of service. In fact, the correlation between

the two variables, `service` and `years` is `r cor(college$service, college$years)`. The variance inflation factor for the five covariates

For categorical variables, the variance inflation factor definition would normally yield for each level a different value; an alternative is the generalized variance inflation factor (Fox & Monette, 1992). Here, we are interested in gender disparities, so the fact that both `service` and `field` are strongly correlated is not problematic, since the VIF for `sex` is not high and the other variables are there to act as control and avoid confounders.

(Generalized) variance inflation factor for the `college` data.

`service`

`years`

`rank`

`sex`

`field`

5.92

7.52

2.01

1.03

1.06

2.9 Graphical analysis of residuals

So far, we have fit models and tested significance of the parameters without checking the model assumptions. The correctness of statements about the p -values and confidence intervals depend on the (approximate) validity of the model assumptions, which all stem from the distributional assumption for the error, assumed to be independent and identically distributed with $\varepsilon_i \sim \text{No}(0, \sigma^2)$. This compact mathematical description can be broken down into four assumptions.

- normality of the errors
- linearity: the mean of Y is $\beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$.
- homoscedasticity: the error variance is constant
- independence of the errors.

This section reviews the assumptions made in order to allow statistical inference using the linear model and different residuals that serve as building blocks for graphical diagnostics.

We investigate the consequences of violation of these assumptions and outline potential mitigation strategies, many of which are undertaken in other chapters.

When we perform an hypothesis test, we merely fail to reject the null hypothesis, either because the latter is true or else due to lack of evidence. The same goes for checking the validity of model assumptions: scientific reasoning dictates that we cannot know for certain whether these hold true. Our strategy is therefore to use implications of the linear model assumptions to create graphical diagnostic tools, so as to ensure that there is no gross violation of these hypothesis. However, it is important to beware of over-interpreting diagnostic plots: the human eye is very good at finding spurious patterns.

Other good references for the material in this section is:

- Forecasting: Principles and Practice, section 5.3

2.9.1 Residuals

Residuals are predictions of the errors ε and represent the difference between the observed value Y_i and the estimated value on the line. The ordinary residuals are

$$e_i = Y_i - \hat{Y}_i, \quad i = 1, \dots, n.$$

The sum of the ordinary residuals is always zero by construction if the model includes an intercept, meaning $\bar{e} = 0$.

Not all observations contribute equally to the adjustment of the fitted hyperplane. The geometry of least squares shows that the residuals are orthogonal to the fitted values, and $e = (\mathbf{I}_n - \mathbf{H}_X)\mathbf{Y}$, where $\mathbf{H}_X = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$ is an $n \times n$ projection matrix that spans the p -dimensional linear combination of the columns of \mathbf{X} , $\mathcal{S}(\mathbf{X})$. If $\text{Va}(\mathbf{Y}) = \sigma^2 \mathbf{I}_n$, it follows that $\text{Va}(e) = \sigma^2(\mathbf{I}_n - \mathbf{H}_X)$ because $(\mathbf{I}_n - \mathbf{H}_X)$ is a projection matrix, therefore idempotent and symmetric. Because the matrix has rank $n - p$, the ordinary residuals cannot be independent from one another.

If the errors are independent and homoscedastic, the ordinary residual e_i has variance $\sigma^2(1 - h_i)$, where the leverage term $h_i = (\mathbf{H}_X)_{ii} = \mathbf{x}_i(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_i$ is the i th diagonal entry of the projection matrix (\mathbf{H}_X) and \mathbf{x}_i is the i th row of the design matrix corresponding to observation i .

We thus conclude that ordinary residuals do not all have the same standard deviation and they are not independent. This is problematic, as we cannot make meaningful comparisons: points with low leverage are bound to deviate more from the fitted model than others. To palliate to this, we can standardize the residuals so each has the same variance under the null of independent homoscedastic errors — the leverage terms h_i are readily calculated from the design matrix \mathbf{X} . The only remaining question is how to estimate the variance. If

we use the i th observation to estimate both the residual and the variance, we introduce additional dependence. A better way is remove the i th observation and refit the model with the $n - 1$ remaining observations to get $s^2_{(-i)}$ (there are tricks and closed-form expressions for these, so one doesn't need to fit n different linear models). The jackknife studentized residual $r_i = e_i / \{s_{(-i)}(1 - h_i)\}$, also termed externally studentized residuals, are not independent, but they are identically distributed and follow a Student distribution with $n - p - 2$ degrees of freedom. These can be obtained in **R** with the command `rstudent`, also in **SAS**.

When to use which residuals? By construction, the vector of ordinary residuals e is orthogonal to the fitted values \hat{y} and also to each column of the design matrix \mathbf{X} : this means a simple linear regression of e with any of these as covariate gives zero intercept and zero slope. However, residual patterns due to forgotten interactions, nonlinear terms, etc. could be picked up from pair plots of ordinary residuals against the explanatories.

While the jackknife studentized residuals r_i are not orthogonal, they are not very different. Use jackknife residuals r to check for equality of variance and distributional assumptions (e.g., using quantile-quantile plots).

One thus typically uses ordinary residuals e for plots of fitted values/explanatories against residuals and otherwise jackknife studentized residuals for any other graphical diagnostic plot.

2.9.2 Leverage and outliers

The leverage h_i of observation i measures its impact on the least square fit, since we can write $h_i = \partial \hat{y}_i / \partial y_i$. Leverage values tell us how much each point impacts the fit: they are strictly positive, are bounded below by $1/n$ and above by 1. The sum of the leverage values is $\sum_{i=1}^n h_i = p + 1$: in a good design, each point has approximately the same contribution, with average weight $(p + 1)/n$.

Points with high leverage are those that have unusual combinations of explanatories. An influential observation ($h_i \approx 1$) pulls the fitted hyperplane towards itself so that $\hat{y}_i \approx y_i$. As a rule of thumb, points with $h_i > 2(p + 1)/n$ should be scrutinized.

It is important to distinguish between **influential** observations (which have unusual x value, i.e., far from the overall mean) and **outliers** (unusual value of the response y). If an observation is both an outlier and has a high leverage, it is problematic.

If influential observations can be detected by inspecting the leverage of each observation, outliers are more difficult to diagnose.

An outlier stands out from the rest of the observations, either because it has an usual response value, or because it falls far from the regression surface. Loosely speaking, an

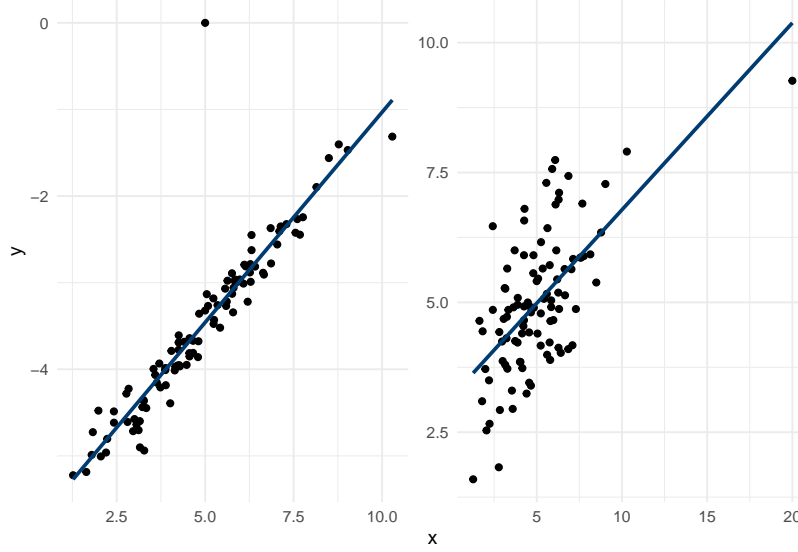


Figure 2.12: Outlier and influential observation. The left panel shows an outlier, whereas the right panel shows an influential variable (rightmost x value).

outlier is an unusual values of Y for a given combination of \mathbf{X} that “stand out” from the rest. Outliers can be detected during the exploratory data analysis or picked-up in residual plots (large values of $|e_i|$ in plots of fitted versus residuals) or added-variable plots. One could potentially test whether a jackknife studentized residual is an outlier (adjusting for the fact we would consider only largest values). One can also consider Cook’s distance, C_j , a statistic giving the scaled distance between the fitted values $\hat{\mathbf{y}}$ and the fitted values for the model with all but the j th observation, $\hat{\mathbf{y}}^{(-j)}$,

$$C_j = \frac{1}{(p+1)S^2} \sum_{i=1}^n \left\{ \hat{y}_i - \hat{y}_i^{(-j)} \right\}^2.$$

Large values of C_j indicate that its residual e_j is large relative to other observations or else its leverage h_j is high. A rule of thumb is to consider points for which $C_j > 4/(n-p-1)$. In practice, if two observations are outlying and lie in the same region, their Cook distance will be halved.

Outliers and influential observations should not be disregarded because they don’t comply with the model, but require further investigation. They may motivate further modelling for features not accounted for. It is also useful to check for registration errors in the data (which can be safely discarded).

Except in obvious scenarios, unusual observations should not be discarded. In very large samples, the impact of a single outlier is hopefully limited. Transformations of the response

may help reduce outlyingness. Otherwise, alternative objective functions (as those employed in robust regression) can be used; these downweight extreme observations, at the cost of efficiency.

2.9.3 Diagnostic plots

We review the assumptions in turn and discuss what happens when the assumptions fail to hold.

2.9.3.1 Independence assumption

Usually, the independence of the observations follows directly from the type of sampling used — this assumption is implicitly true if the observations were taken from a *random sample* from the population. This is generally not the case for longitudinal data, which contains repeated measures from the same individuals across time. Likewise, time series are bound not to have independent observations. If we want to include all the time points in the analysis, we must take into account the possible dependence (correlation) between observations. If we ignore correlation, the estimated standard errors are too small relative to the truth, so the effective sample size is smaller than number of observations.

The first source of dependence is clustered data, meaning measurements taken from subjects that are not independent from one another (family, groups, etc.) More generally, correlation between observations arises from time dependence, roughly categorized into

- longitudinal data: repeated measurements are taken from the same subjects (few time points)
- time series: observations observed at multiple time periods (many time points). Time series require dedicated models not covered in this course.

Because of autocorrelation, positive errors tend to be followed by positive errors, etc. We can plot the residuals as a function of time, and a scatterplot of lagged residuals e_i versus e_{i-1} ($i = 2, \dots, n$).

However, lagged residuals plots only show dependence at lag one between observations. For time series, we can look instead at a correlogram, i.e., a bar plot of the correlation between two observations h units apart as a function of the lag h (Brockwell & Davis, 2016, Definition 1.4.4).

For y_1, \dots, y_n and constant time lags $h = 0, 1, \dots$ units, the autocorrelation at lag h is

$$r(h) = \frac{\gamma(h)}{\gamma(0)}, \quad \gamma(h) = \frac{1}{n} \sum_{i=1}^{n-|h|} (y_i - \bar{y})(y_{i+h} - \bar{y})$$

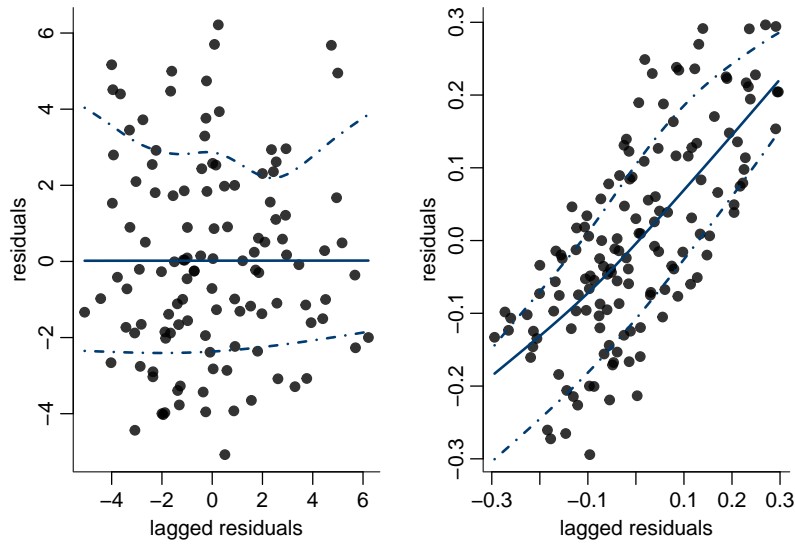


Figure 2.13: Lagged residual plots: there is no evidence against independence in the left panel, whereas the right panel shows positively correlated residuals.

If the series is correlated, the sample autocorrelation will likely fall outside of the pointwise confidence intervals, as shown in Figure 2.14. Presence of autocorrelation requires modelling the correlation between observations explicitly using dedicated tools from the time series literature that are covered in MATH 60638. We will however examine AR(1) models as part of the chapter on longitudinal data.

When observations are positively correlated, the estimated standard errors reported by the software are too small. This means we are overconfident and will reject the null hypothesis more often than we should if the null is true (inflated Type I error, or false positive).

2.9.3.2 Linearity assumption

The second assumption of the linear model is that of linearity, which means that the mean model is correctly specified, all relevant covariates have been included and their effect is correctly specified. To check that the response surface of the linear model is adequate, we plot e_i against \hat{y}_i or X_{ij} (for $j = 1, \dots, p$). Since the linear correlation between e and \hat{y} (or e and X_j) is zero by construction, patterns (e.g., quadratic trend, cycles, changepoints) are indicative of misspecification of the mean model. One can add a smoother to detect patterns. Figure 2.15 shows three diagnostics plots, the second of which shows no pattern in the residuals, but skewed fitted values.

If there is residual structure in plots of ordinary residuals against either (a) the fitted values or

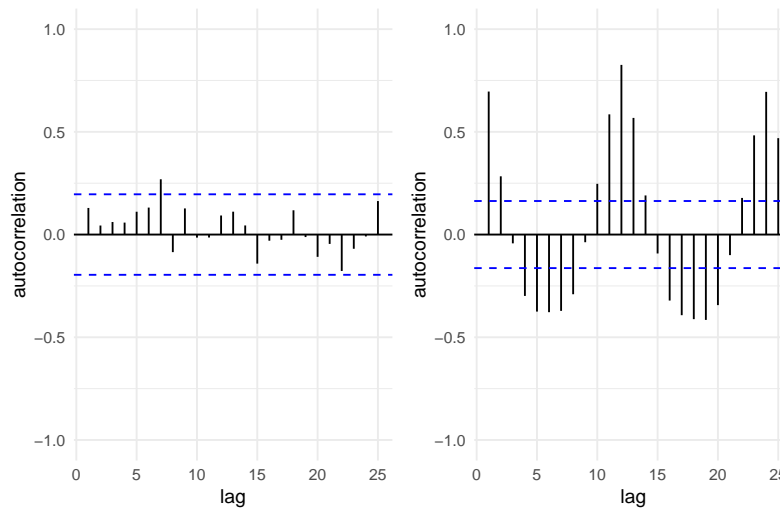


Figure 2.14: Correlogram of independent observations (left) and the ordinary residuals of the log-linear model fitted to the air passengers data (right). While the mean model of the latter is seemingly correctly specified, there is residual dependence between monthly observations and yearly (at lag 12). The blue lines give approximate pointwise 95% confidence intervals for white noise (uncorrelated observations).

(b) the explanatory variables, a more complex model can be adjusted including interactions, nonlinear functions, ... If the effect of an explanatory variable is clearly nonlinear and complicated, smooth terms could be added (we won't cover generalized additive models in this course).

Plotting residuals against left-out explanatory variables can also serve to check that all of the explanatory power of the omitted covariate is already explained by the columns of \mathbf{X} .

If an important variable has been omitted and is not available in the dataset, then the effect of that variable is captured by both the errors (the portion orthogonal to the design matrix \mathbf{X} , i.e., unexplained by the covariates included in the model) and the remaining part is captured by other explanatory variables of the model that are correlated with the omitted variable. These variables can act as confounders. There is little that can be done in either case unless the data for the omitted variable are available, but subject-specific knowledge may help make sense of the results.

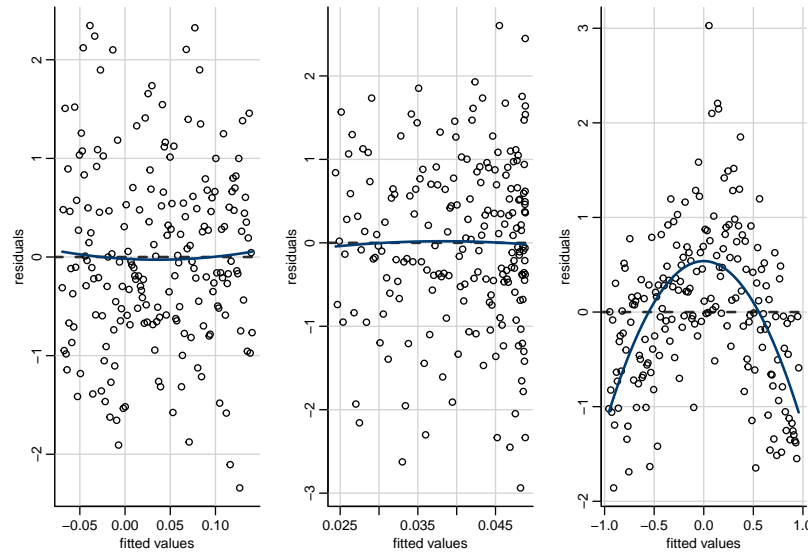


Figure 2.15: Scatterplots of residuals against fitted values. The first two plots show no departure from linearity (mean zero). The third plot shows a clear quadratic pattern, suggesting the mean model is misspecified. Note that the distribution of the fitted value need not be uniform, as in the second panel which shows more high fitted values.

2.9.3.3 Homoscedasticity assumption

If the variance of the errors is the same for all observations, that of the observations Y is also constant. The most common scenarios for heteroscedasticity are increases in variance with the response, or else variance that depends on explanatory variables \mathbf{X} , most notably categorical variables. For the former, a log-transform (or Box–Cox transformation) can help stabilize the variance, but we need the response to be positive. For the latter, we can explicitly model that variance and we will see how to include different variance per group later on. A popular strategy in the econometrics literature, is to use robust (inflated) estimators of the standard errors such as White’s sandwich estimator of the variance.

If the residuals (or observations) are heteroscedastic (non constant variance), the estimated effects of the variables (the β parameters) are still valid in the sense that the ordinary least squares estimator $\hat{\beta}$ is unbiased. However, the estimated standard errors of the $\hat{\beta}$ are no longer reliable and, consequently, the confidence intervals and the hypothesis tests for the model parameters will be incorrect. Indeed, if the variance of the errors differs from one observation to the next, we will estimate an average of the different variance terms. The standard errors of each term are incorrect (too small or too large) and the conclusions of the tests (p -values) will be off because the formulas of both t -test and F -test statistics include

estimates of $\hat{\sigma}^2$.

Looking at the plot of jackknife studentized residuals against regressors (or fitted values) is instructive — for example, we often see a funnel pattern when there is an increase in variance in the plot of the jackknife studentized residuals against fitted value, or else in boxplots with a categorical variable as in Figure 2.17. However, if we want to fit a local smoother to observe trends, it is better to plot the absolute value of the jackknife studentized residuals against regressors or observation number.

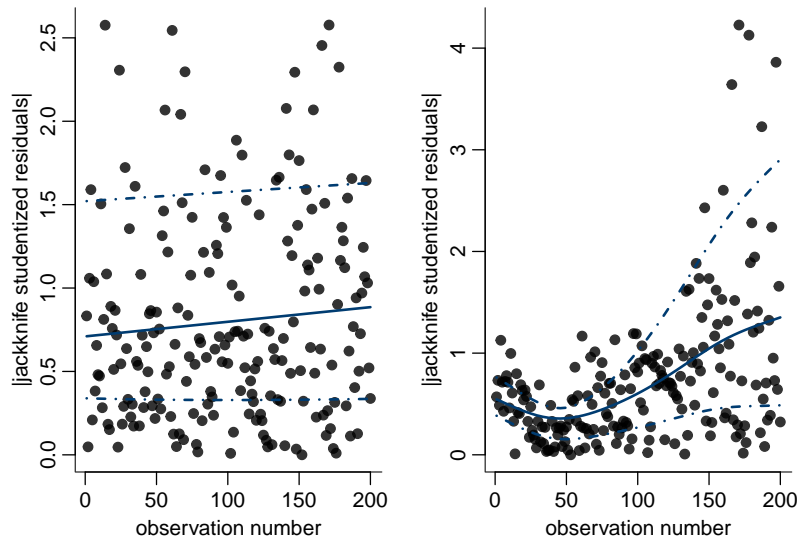


Figure 2.16: Plot of the absolute value of jackknife studentized residuals against observation number. The left panel is typical of homoscedastic data, whereas the right panel indicates an increase in the variance.

2.9.3.4 Normality assumption

The normality assumption is mostly for convenience: if the errors are assumed normally distributed, then the least square and the maximum likelihood estimators of β coincide. The maximum likelihood estimators of β are asymptotically normal under mild conditions on the design matrix and t -tests are robust to departure of the normality assumption. This means that inference is valid in large samples, regardless of the distribution of the errors/residuals (even if the null distribution are not exact). It is important to keep in mind that, for categorical explanatory variables, the sample size in each group must be sufficiently large for the central limit theorem to kick in.

Sometimes, transformations can improve normality: if the data is right-skewed and the

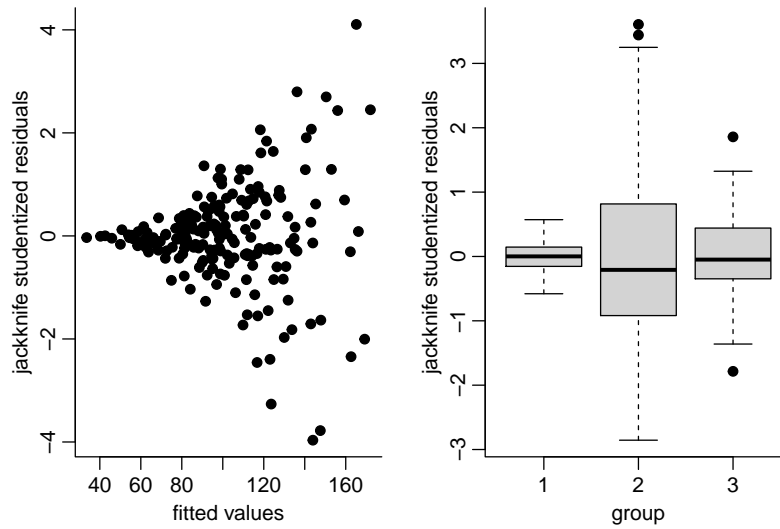


Figure 2.17: Plot of jackknife studentized residuals against fitted value (left) and categorical explanatory (right). Both clearly display heteroscedasticity.

response is strictly positive, a log-linear model may be more adequate. This can be assessed by looking at the quantile-quantile plot of the externally studentized residuals. If the response Y is not continuous (including binary, proportion or count data), linear models give misleading answers and generalized linear models are more suitable.

The inference will be valid for large samples even if the errors are not normally distributed by virtue of the central limit theorem. If the errors $\varepsilon_i \sim \text{No}(0, \sigma^2)$, then the jackknife studentized residuals should follow a Student distribution, with $r_i \sim \text{St}(n-p-2)$, (identically distributed, but not independent). A Student quantile-quantile plot can thus be used to check the assumption (and for n large, the normal plotting positions could be used as approximation if $n-p > 50$). One can also plot a histogram of the residuals. Keep in mind that if the mean model is not correctly specified, some residuals may incorporate effect of leftover covariates.

Quantile-quantile plots are discussed in Section A.2.5, but their interpretation requires training. For example, Figure 2.19 shows many common scenarios that can be diagnosed using quantile-quantile plots: discrete data is responsible for staircase patterns, positively skewed data has too high low quantiles and too low high quantiles relative to the plotting positions, heavy tailed data have high observations in either tails and bimodal data leads to jumps in the plot.

Example 2.4 (Diagnostic plots for the college data.). We can look at the college data to

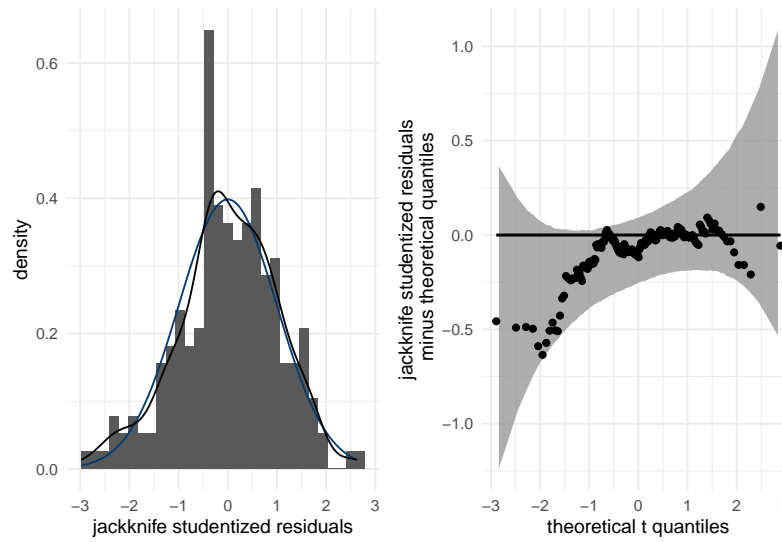


Figure 2.18: Histogram (left) and Student quantile-quantile plot (right) of the jackknife studentized residuals. The left panel includes a kernel density estimate (black), with the density of Student distribution (blue) superimposed. The right panel includes pointwise 95% confidence bands calculated using a bootstrap.

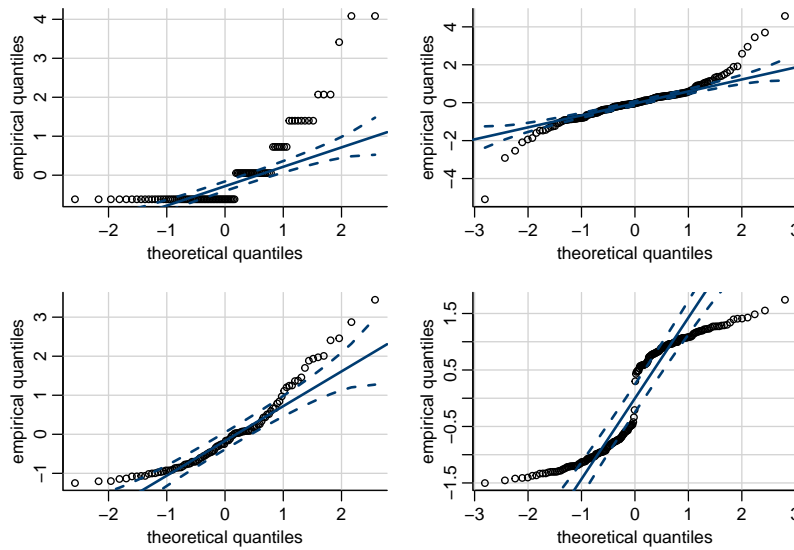


Figure 2.19: Quantile-quantile plots of non-normal data, showing typical look of behaviour of discrete (top left), heavy tailed (top right), skewed (bottom left) and bimodal data (bottom right).

see if the linear model assumptions hold.

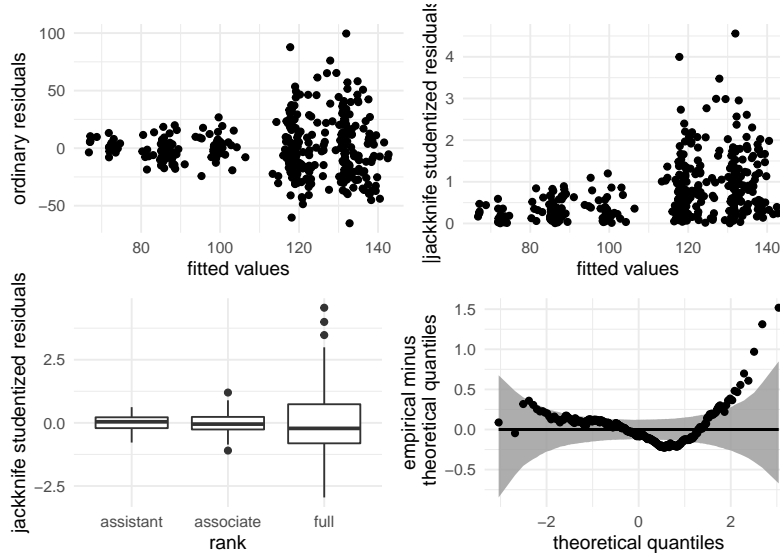


Figure 2.20: Diagnostic plots for the college data example: ordinary residuals against fitted values (top left), absolute value of the jackknife studentized residuals against fitted values (top right), box and whiskers plot of jackknife studentized residuals (bottom left) and detrended Student quantile-quantile plot (bottom right). There is clear group heteroscedasticity.

Based on the plots of Figure 2.20, we find that there is residual heteroscedasticity, due to rank. Since the number of years in the first rank is limited and all assistant professors were hired in the last six years, there is less disparity in their income. It is important not to mistake the pattern on the x -axis for the fitted value (due to the large effect of rank and field, both categorical variable) with patterns in the residuals (none apparent). Fixing the heteroscedasticity would correct the residuals and improve the appearance of the quantile-quantile plot.

2.9.4 Model adjustments

This ultimate section deals with strategies for fixing the linear model if we detected non-normality. If the response is strictly positive, an option is to use the Box–Cox transformation presented in Section 3.3 and fit a linear model to a transformation of Y ,

$$y_i(\lambda) = \begin{cases} (y^\lambda - 1)/\lambda, & \lambda \neq 0 \\ \ln(y), & \lambda = 0. \end{cases}$$

We then fit a linear model over a grid of λ , using the profile likelihood to select the optimal value for the transformation. Note that the latter will depend on the covariates present in the model; the example with the `college` data already shows that the diagnostics of normality is impacted by departures from the other hypothesis. The cases $\lambda = 1$ (identity) and $\lambda = 0$ (log-linear model) are perhaps the most important because they yield interpretable coefficients.

If the data is right-skewed and the response is strictly positive, a log-linear model may be more adequate and the parameters can be interpreted. Theory sometimes dictates a multiplicative model: for example, the Cobb–Douglas production function in economics is $P = \alpha L^{\beta_1} C^{\beta_2}$, where P stands for production, L for labor and C for capital; all inputs are positive, so taking a log-transform yields a model that is linear in β , with $\beta_0 = \ln(\alpha)$.

We can rewrite the model in the original response scale as

$$Y = \exp\left(\beta_0 + \sum_{j=1}^p \beta_j X_j + \varepsilon\right) = \exp\left(\beta_0 + \sum_{j=1}^p \beta_j X_j\right) \cdot \exp(\varepsilon),$$

and thus

$$E(Y \mid \mathbf{X}) = \exp(\beta_0 + \beta_1 X_1 + \cdots + \beta_p X_p) \times E\{\exp(\varepsilon) \mid \mathbf{X}\}.$$

If $\varepsilon \mid \mathbf{X} \sim \text{No}(\mu, \sigma^2)$, then $E\{\exp(\varepsilon) \mid \mathbf{X}\} = \exp(\mu + \sigma^2/2) = \exp(\varepsilon)$ follows a log-normal distribution.

In order to interpret the parameters, we can compare the ratio of $E(Y \mid X_1 = x + 1)$ to $E(Y \mid X_1 = x)$,

$$\frac{E(Y \mid X_1 = x + 1, X_2, \dots, X_p)}{E(Y \mid X_1 = x, X_2, \dots, X_p)} = \frac{\exp\{\beta_1(x + 1)\}}{\exp(\beta_1 x)} = \exp(\beta_1).$$

Thus, $\exp(\beta_1)$ represents the ratio of the mean of Y when $X_1 = x + 1$ in comparison to that when $X_1 = x$, *ceteris paribus* (and provided this statement is meaningful). We can interpret $\exp(\beta_1)$ as the multiplicative effect of X_1 on the mean of Y : increasing X_1 by one unit causes Y to increase by a factor of $\exp(\beta_1)$, on average.

Example 2.5. The paper of Box–Cox consider survival time for 48 animals based on a randomized trial and these are analyzed in Example 8.25 of Davison (2008). Three poisons were administered with four treatments; each factor combination contained four animals, chosen at random. There is strong evidence that both the choice of poison and treatment affect survival time.

We could consider a two-way analysis of variance model for these data without interaction, given the few observations for each combination. The model would be of the form

$$Y = \beta_0 + \beta_1 \text{poison}_2 + \beta_2 \text{poison}_3 + \beta_3 \text{treatment}_2 \\ + \beta_4 \text{treatment}_3 + \beta_5 \text{treatment}_4 + \varepsilon$$

The plot of fitted values against residuals shows that the model is not additive; there is also indications that the variance increases with the mean response. The model is inadequate: lowest survival times are underpredicted, meaning the residuals are positive and likewise the middle responses is positive. A formal test of non-additivity based on constructed variables further point towards non-additivity (Davison, 2008, Example 8.24). All of these factors point towards overall poor fit of the model.

One could ignore the non-constant variance and consider using a Box–Cox to find a suitable transformation of the residuals so as to improve normality. The profile log likelihood at the bottom left of Figure 2.21 suggests that $\lambda \approx -1$ would be a good choice. This has the benefit of being interpretable, as the reciprocal response Y^{-1} corresponds to the speed of action of the poison depending on both poison type and treatment. The diagnostics plots also indicate that the model for the reciprocal has no residual structure and the variance appears constant.

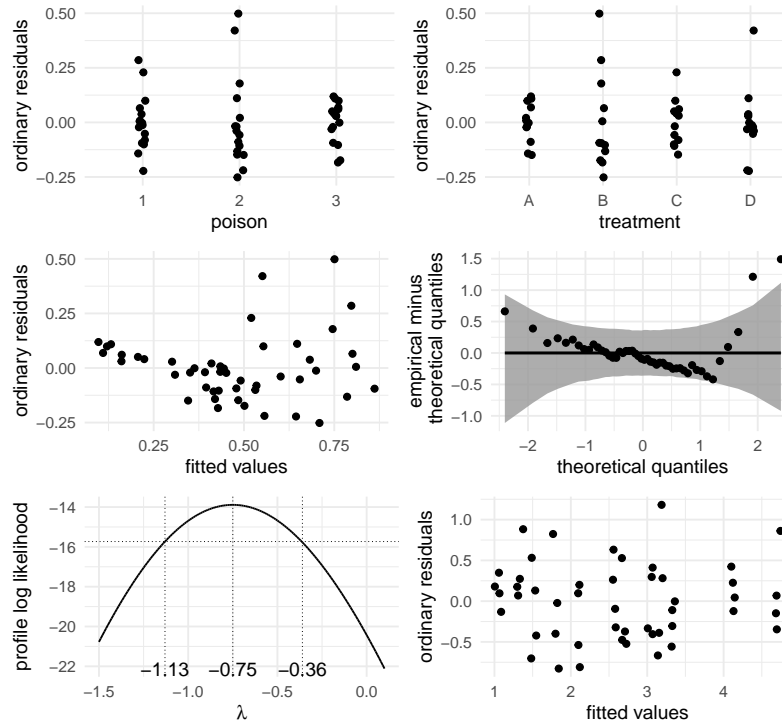


Figure 2.21: Diagnostic plots for the poison data. The top panel shows the ordinary residuals for the linear model for survival time as a function of poison and treatment, with jittered observations. The middle left plot shows the fitted values against residuals, which display evidence of trend and increase in variance with the survival time. The quantile-quantile plot in the middle right plot shows some evidence of departure from the normality, but the non-linearity and heteroscedasticity obscure this. The bottom panel shows the profile log likelihood for the Box–Cox transform, suggesting a value of -1 would be within the 95% confidence interval. After fitting the same additive model with main effect only to the reciprocal survival time, there is no more evidence of residual structure and unequal variance.

Chapter 3

Likelihood-based inference

The goal of this chapter is to familiarize you with likelihood-based inference.

The starting point of likelihood-based inference is a statistical model: we postulate that (a function of) the data has been generated from a probability distribution with p -dimensional parameter vector θ . The purpose of the analyst is to estimate these unknown parameters on the basis of a sample and make inference about them.

3.1 Maximum likelihood

The **likelihood** $L(\theta)$ is a function of θ that gives the probability (or density) of observing a sample under a postulated distribution, treating the observations as fixed. In most settings we consider, observations are independent and so the joint probability of the sample values is the product of the probability of the individual observations¹: for y_1, \dots, y_n assuming Y_i ($i = 1, \dots, n$) follows a distribution whose mass function or density is $f(y; \theta)$, this is just

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n f(y_i; \theta) = f(y_1; \theta) \times \dots \times f(y_n; \theta).$$

From a pure optimization perspective, the likelihood is a particular choice of objective function that reflects the probability of the observed outcome. One shouldn't however maximize directly the likelihood, since computing the product of a lot of potentially small numbers is subject to numerical overflow and is unstable (for discrete distributions, the mass function gives probabilities that are by definition between zero and one). Instead, one

¹If this seems foreign, think about repeated coin tosses with Bernoulli distribution of unknown parameter p and convince yourself that the trials are independent, so the probability of obtaining two consecutive heads is 0.25 for a fair coin.

should work with the log likelihood function, $\ell(\theta) = \ln\{L(\theta)\}$. Since logarithm is a strictly increasing function, maximizing the natural logarithm (denoted \ln) of the likelihood leads to the same solution. Another reason why working with the log likelihood is preferable is because product over n likelihood contributions becomes a sum and this facilitates numerical and analytical derivations of the maximum likelihood estimators (the log of a product is equal to the sum of logs, i.e., $\ln(ab) = \ln(a) + \ln(b)$ for $a, b > 0$.)

The **maximum likelihood estimator** $\hat{\theta}$ is the value of θ that maximizes the likelihood, i.e., the value under which the random sample is the most likely to be generated. The scientific reasoning behind this is: “whatever we observe, we have expected to observe” so we choose between competing models the one that makes the most sense.

Several properties of maximum likelihood estimator makes it appealing for inference.

- The maximum likelihood estimator is **consistent**, i.e., it converges to the correct value as the sample size increase (asymptotically unbiased).
- The maximum likelihood estimator is invariant to reparametrizations
- Under regularity conditions, the maximum likelihood estimator is asymptotically normal, so we can obtain the null distribution of classes of hypothesis tests and derive confidence intervals based on $\hat{\theta}$.
- The maximum likelihood estimator is efficient, meaning it has the smallest asymptotic mean squared error (or the smallest asymptotic variance).

The **score function** $U(\theta; \mathbf{y}) = \partial\ell(\theta; \mathbf{y})/\partial\theta$ is the gradient of the log likelihood function and, under regularity conditions, the maximum likelihood estimator solves $U(\theta; \mathbf{Y}) = \mathbf{0}_p$. This property can be used to derive gradient-based algorithms for optimization and for verifying that the solution found is a global maximum.

Remark. While least squares admit a closed-form solution, the maximum of the log likelihood is generally found numerically by solving the score equation. The algorithms used in most software are reliable and efficient for regression models we consider in this course. However, for more complex models, like generalized linear mixed models, the convergence of optimization algorithms is oftentimes problematic and scrutiny is warranted.

The **observed information matrix** is the hessian $j(\theta; \mathbf{y}) = -\partial^2\ell(\theta; \mathbf{y})/\partial\theta\partial\theta^\top$ evaluated at the maximum likelihood estimate $\hat{\theta}$. Under regularity conditions, the Fisher information matrix is

$$i(\theta) = E\{U(\theta; \mathbf{Y})U(\theta; \mathbf{Y})^\top\} = E\{j(\theta; \mathbf{Y})\}$$

The Fisher (or expected) and observed information matrices encodes the curvature of the log likelihood and provides information about the variability of $\hat{\theta}$.

The properties of the log likelihood are particularly convenient for inference because they provide omnibus testing procedures that have a known asymptotic distribution. The

starting point for the distributional theory surrounding likelihood-based statistics is the asymptotic normality of the score $U(\boldsymbol{\theta}) \sim \text{No}(0, i(\boldsymbol{\theta}))$, which follows from a central limit theorem. The variance of $U(\boldsymbol{\theta}_0)$ is exactly $i(\boldsymbol{\theta}_0)$, while that of $\hat{\boldsymbol{\theta}}$ is approximately $i(\boldsymbol{\theta}_0)^{-1}$ under the null hypothesis \mathcal{H}_0 . This result is particularly useful: we often use the inverse of the observed information as estimate of the covariance matrix of the maximum likelihood estimator. To obtain the standard errors of $\hat{\boldsymbol{\theta}}$, one simply computes the square root of the diagonal elements of the inverse of the observed information, i.e., $[\text{diag}\{j^{-1}(\hat{\boldsymbol{\theta}})\}]^{1/2}$.

Example 3.1 (Exponential model for waiting times of the Montreal metro). Consider the waiting time Y between consecutive subways arriving at Station Édouard-Montpetit on the blue line in Montreal during rush hour. We postulate that these waiting times follow an exponential distribution with rate θ , denoted $Y \sim E(\theta)$. The purpose of statistical inference is to use the information from a random sample of size n to estimate the unknown parameter θ . The density of Y evaluated at y , $f(y; \theta) = \theta \exp(-\theta y)$, encodes the probability of the observed waiting time for a given parameter value and, if the records are independent, the probability of observing y_1, \dots, y_n is the product of probabilities of individual events. The likelihood is thus

$$L(\theta; \mathbf{y}) = \prod_{i=1}^n f(y_i; \theta) = \prod_{i=1}^n \theta \exp(-\theta y_i),$$

$$\ell(\theta; \mathbf{y}) = n \ln(\theta) - \theta \sum_{i=1}^n y_i$$

To find the maximum of the function, we differentiate the log likelihood $\ell(\theta; \mathbf{y})$ and set the gradient to zero,

$$\frac{\partial \ell(\theta; \mathbf{y})}{\partial \theta} = \frac{n}{\theta} - \sum_{i=1}^n y_i = 0.$$

Solving for θ gives $\hat{\theta} = \bar{y}^{-1}$, so the maximum likelihood estimator is the reciprocal of the sample mean \bar{Y} . The observed information is $j(\theta) = n\theta^{-2}$ and likewise $i(\theta) = E\{j(\theta)\} = n\theta^{-2}$.

For the sample of waiting time in the subway, the maximum likelihood estimate of θ is $\hat{\theta} = 0.327$, the observed information is $j(\hat{\theta}) = i(\hat{\theta}) = 0.467$ and the standard error of $\hat{\theta}$ is $j(\hat{\theta})^{-1/2} = 1.463$.

Example 3.2 (Normal samples and ordinary least squares). Suppose we have an independent normal sample of size n with mean μ and variance σ^2 , where $Y_i \sim \text{No}(\mu, \sigma^2)$ are independent. Recall that the density of the normal distribution is

$$f(\mathbf{y}; \mu, \sigma^2) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (x - \mu)^2 \right\}.$$

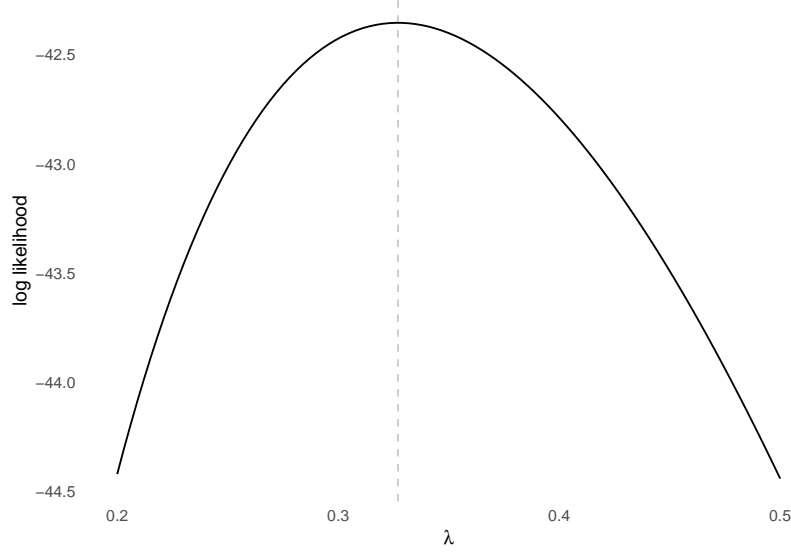


Figure 3.1: Log-likelihood for a sample of size 20 of waiting times (in minutes)

For an n -sample \mathbf{y} , the likelihood is

$$\begin{aligned} L(\mu, \sigma^2; \mathbf{y}) &= \prod_{i=1}^n \frac{1}{(2\pi\sigma^2)^{1/2}} \exp \left\{ -\frac{1}{2\sigma^2} (y_i - \mu)^2 \right\} \\ &= (2\pi\sigma^2)^{-n/2} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2 \right\}. \end{aligned}$$

and the log likelihood is

$$\ell(\mu, \sigma^2; \mathbf{y}) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - \mu)^2.$$

One can show that the maximum likelihood estimators for the two parameters are

$$\hat{\mu} = \bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

The fact that the estimator of the theoretical mean μ is the sample mean is fairly intuitive and one can show the estimator is unbiased for μ . The (unbiased) sample variance estimator,

$$S^2 = \frac{1}{n-1} \sum_{i=1}^n (Y_i - \bar{Y})^2$$

Since $\hat{\sigma}^2 = (n-1)/nS^2$, it follows that the maximum likelihood estimator of σ^2 is biased, but both estimators are consistent and will thus get arbitrarily close to the true value σ^2 for n sufficiently large.

The case of normally distributed data is intimately related to linear regression and ordinary least squares: assuming normality of the errors, the least square estimators of β coincide with the maximum likelihood estimator of β .

Recall the linear regression model,

$$Y_i = \beta_0 + \beta_1 X_{i1} + \beta_2 X_{i2} + \dots + \beta_p X_{ip} + \varepsilon_i, \quad (i = 1, \dots, n),$$

where the errors $\varepsilon_i \sim \text{No}(0, \sigma^2)$. The linear model has $p+2$ parameters (β and σ^2) and the log likelihood is

$$\ell(\theta) = -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \{ \mathbf{y} - \mathbf{X}\beta \}^\top (\mathbf{y} - \mathbf{X}\beta) \}^2.$$

Maximizing the log likelihood with respect to β is equivalent to minimizing the sum of squared errors $\|\mathbf{Y} - \widehat{\mathbf{Y}}\|^2$. Since this objective function is the same as that of least squares, it follows that the least-square estimator $\hat{\beta}$ for the mean parameters is the maximum likelihood estimator for normal errors with common variance σ^2 , regardless of the value of the latter. The maximum likelihood estimator $\hat{\sigma}^2$ is thus

$$\hat{\sigma}^2 = \max_{\sigma^2} \ell(\hat{\beta}, \sigma^2).$$

The log likelihood, excluding constant terms that don't depend on σ^2 , is

$$\ell(\hat{\beta}, \sigma^2) \propto -\frac{1}{2} \left\{ n \ln \sigma^2 + \frac{1}{\sigma^2} (\mathbf{y} - \mathbf{X}\hat{\beta})^\top (\mathbf{y} - \mathbf{X}\hat{\beta}) \right\}.$$

Differentiating each term with respect to σ^2 and setting the gradient equal to zero yields the maximum likelihood estimator

$$\hat{\sigma}^2 = \frac{1}{n} (\mathbf{Y} - \mathbf{X}\hat{\beta})^\top (\mathbf{Y} - \mathbf{X}\hat{\beta}) = \frac{1}{n} \sum_{i=1}^n e_i^2 = \frac{\text{SS}_e}{n};$$

where SS_e is the sum of squared residuals. The usual unbiased estimator of σ^2 calculated by software is $S^2 = \text{SS}_e/(n-p-1)$, where the denominator is the sample size n minus the number of mean parameters β , $p+1$.

Sometimes, the p -parameter vector θ of the likelihood is not the quantity of interest. Suppose for simplicity we are interested in a scalar function $\phi = g(\theta)$. The maximum likelihood estimate of ϕ is $\hat{\phi} = g(\hat{\theta})$. This property of maximum likelihood estimators justifies their widespread use. In large samples, $\hat{\theta}$ is centered at the true value θ_0 and is approximately

multivariate normal with $\hat{\boldsymbol{\theta}} \sim \text{No}_p(\boldsymbol{\theta}_0, \mathbf{V}_{\boldsymbol{\theta}})$, then $\hat{\phi} \sim \text{No}(\phi_0, V_{\phi})$, with $V_{\phi} = \nabla \phi^{\top} \mathbf{V}_{\boldsymbol{\theta}} \nabla \phi$, where $\nabla \phi = [\partial \phi / \partial \theta_1, \dots, \partial \phi / \partial \theta_p]^{\top}$. In applications, the variance matrix and the gradient vector are evaluated at the maximum likelihood estimate $\hat{\boldsymbol{\theta}}$

Consider the metro waiting time example. The quantity of interest is the reciprocal mean $\phi = 1/\theta$, so the scalar function of interest is $g(x) = 1/x$ and the maximum likelihood estimate $\hat{\phi} = 3.058$ is the sample mean of waiting times. The gradient of the transformation is $\nabla \phi = -1/x^2$, which gives $V_{\phi} = (\theta^2/n)/\theta^4 = 1/(n\theta^2) = \phi^2/n$ and the standard error of the maximum likelihood estimator is thus approximately $\hat{\phi}/\sqrt{n}$ in large samples.

3.2 Likelihood-based tests

Oftentimes, we wish to compare two models: the model implied by the null hypothesis, which is a restriction or simpler version of the full model. Models are said to be **nested** if we can obtain one from the other by imposing restrictions on the parameters.

We consider a null hypothesis \mathcal{H}_0 that imposes restrictions on the possible values of $\boldsymbol{\theta}$ can take, relative to an unconstrained alternative \mathcal{H}_1 . We need two **nested** models: a *full* model, and a *reduced* model that is a subset of the full model where we impose q restrictions. For example, the full model could be a regression model with four predictor variables and the reduced model could include only the first two predictor variables, which is equivalent to setting $\mathcal{H}_0 : \beta_3 = \beta_4 = 0$. The testing procedure involves fitting the two models and obtaining the maximum likelihood estimators of each of \mathcal{H}_1 and \mathcal{H}_0 , respectively $\hat{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}_0$ for the parameters under \mathcal{H}_0 . The null hypothesis \mathcal{H}_0 tested is: ‘the reduced model is an **adequate simplification** of the full model’ and the likelihood provides three main classes of statistics for testing this hypothesis: these are

- likelihood ratio tests statistics, denoted R , which measure the drop in log likelihood (vertical distance) from $\ell(\hat{\boldsymbol{\theta}})$ and $\ell(\hat{\boldsymbol{\theta}}_0)$.
- Wald tests statistics, denoted W , which consider the standardized horizontal distance between $\hat{\boldsymbol{\theta}}$ and $\hat{\boldsymbol{\theta}}_0$.
- score tests statistics, denoted S , which looks at the scaled gradient of ℓ , evaluated *only* at $\hat{\boldsymbol{\theta}}_0$ (derivative of ℓ).

The three main classes of statistics for testing a simple null hypothesis $\mathcal{H}_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$ against the alternative $\mathcal{H}_a : \boldsymbol{\theta} \neq \boldsymbol{\theta}_0$ are the likelihood ratio, the score and the Wald statistics, defined respectively as

$$\begin{aligned} R &= 2 \left\{ \ell(\hat{\boldsymbol{\theta}}) - \ell(\boldsymbol{\theta}_0) \right\}, \\ S &= U^{\top}(\boldsymbol{\theta}_0) i^{-1}(\boldsymbol{\theta}_0) U(\boldsymbol{\theta}_0), \\ W &= (\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0)^{\top} j(\hat{\boldsymbol{\theta}})(\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}_0), \end{aligned}$$

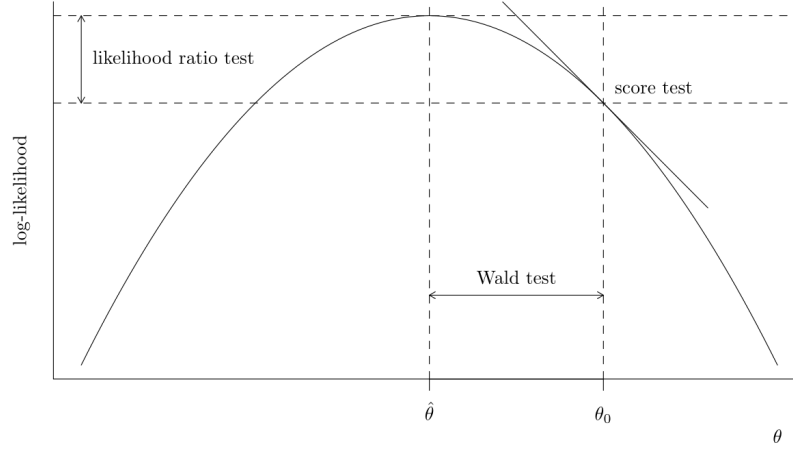


Figure 3.2: Log-likelihood curve: the three likelihood-based tests, namely Wald, likelihood ratio and score tests, are shown on the curve. The tests use different information about the function.

where $\hat{\theta}$ is the maximum likelihood estimate under the alternative and θ_0 is the null value of the parameter vector. Asymptotically, all the test statistics are equivalent (in the sense that they lead to the same conclusions about \mathcal{H}_0). If \mathcal{H}_0 is true, the three test statistics follow asymptotically a χ_q^2 distribution under a null hypothesis \mathcal{H}_0 , where the degrees of freedom q are the number of restrictions.

For scalar θ with $q = 1$, signed versions of these statistics exist, e.g.,

$$W(\theta_0) = (\hat{\theta} - \theta_0)/\text{se}(\hat{\theta}) \sim \text{No}(0, 1)$$

for the Wald statistic or the directed likelihood root

$$R(\theta_0) = \text{sign}(\hat{\theta} - \theta_0) \left[2 \{ \ell(\hat{\theta}) - \ell(\theta_0) \} \right]^{1/2} \sim \text{No}(0, 1).$$

The likelihood ratio test statistic is normally the most powerful of the three likelihood tests. The score statistic S only requires calculation of the score and information under \mathcal{H}_0 (because by definition $U(\hat{\theta}) = 0$), so it can be useful in problems where calculations of the maximum likelihood estimator under the alternative is costly or impossible.

The Wald statistic W is the most widely encountered statistic and two-sided 95% confidence intervals for a single parameter θ are of the form

$$\hat{\theta} \pm q_{1-\alpha/2} \text{se}(\hat{\theta}),$$

where $q_{1-\alpha/2}$ is the $1-\alpha/2$ quantile of the standard normal distribution; for a 95% confidence interval, the 0.975 quantile of the normal distribution is 1.96. The Wald-based confidence

intervals are by construction **symmetric**: they may include implausible values (e.g., negative values for variances). The Wald-based confidence intervals are not parametrization invariant: if we want intervals for a nonlinear continuous function $g(\theta)$, then in general $\text{Cl}_W\{g(\theta)\} \neq g\{\text{Cl}_W(\theta)\}$.

These confidence intervals can be contrasted with the (better) ones derived using the likelihood ratio test: these are found through a numerical search to find the limits of

$$\theta : 2\{\ell(\hat{\theta}) - \ell(\theta)\} \leq \chi_1^2(1 - \alpha),$$

where $\chi_1^2(1 - \alpha)$ is the $(1 - \alpha)$ quantile of the χ_1^2 distribution. If θ is multidimensional, confidence intervals for θ_i are derived using the profile likelihood. Likelihood ratio-based confidence intervals are **parametrization invariant**, so $\text{Cl}_R\{g(\theta)\} = g\{\text{Cl}_R(\theta)\}$. Because the likelihood is zero if a parameter value falls outside the range of possible values for the parameter, the intervals only include plausible values of θ . In general, the intervals are asymmetric and have better coverage properties.

To illustrate the difference between likelihood ratio and Wald tests (and their respective confidence intervals), we consider the metro waiting time data and consider a two-sided test for the null hypothesis $\mathcal{H}_0 : \phi_0 = 4$ minutes. The Wald statistic is

$$W = \frac{\hat{\phi} - \phi_0}{\text{se}(\hat{\phi})} = \sqrt{n} \frac{\hat{\phi} - \phi_0}{\hat{\phi}} = -1.378,$$

since $\text{se}(\phi) = V_\phi^{1/2} = \phi/\sqrt{n}$ and the latter function is evaluated at $\hat{\phi} = 3.058$. The asymptotic null distribution is $\text{No}(0, 1)$, so we fail to reject $\mathcal{H}_0 : \phi = 4$ minutes since the observed value for $|W|$ is smaller than 1.96. We could invert the test statistic to get a symmetric 95% confidence interval for ϕ , $[1.718, 4.398]$.

The hypothesis corresponds also to $\theta = 0.25$ and similar calculations give $W = \sqrt{n}(0.327 - 0.25)/0.327 = 1.053$. In this case, the null distribution is the same, but the value of the test statistic is not! The confidence interval for θ is $[0.184, 0.47]$. You can check for yourself that the reciprocal of these confidence intervals do not match those for ϕ .

In contrast, the likelihood ratio test is invariant to interest-preserving reparametrizations, so the test statistic for $\mathcal{H}_0 : \phi = 4$ and $\mathcal{H}_0 : \theta = 0.25$ are the same. The log likelihood at $\theta = 0.25$ is -43.015, the maximum log likelihood value is -42.354 and the likelihood ratio statistic $R = 2\{\ell(\hat{\theta}) - \ell(4)\} = 1.322$. The statistic should behave like a χ_1^2 variable in large samples. The 95% of the χ_1^2 distribution is 3.841, so we fail to reject the null hypothesis that the mean waiting time is 4 minutes.

The likelihood ratio statistic 95% confidence interval for ϕ can be found by using a root finding algorithm: the confidence interval is $[2.032, 4.906]$. By invariance, the 95% confidence interval for θ is $[0.204, 0.492]$.

3.3 Profile likelihood

Sometimes, we may want to perform hypothesis test or derive confidence intervals for selected components of the model. For example, we may be interested in obtaining confidence intervals for a single β_j in a logistic regression, treating the other parameters β_{-j} as nuisance. In this case, the null hypothesis only restricts part of the space and the other parameters, termed nuisance, are left unspecified — the question then is what values to use for comparison with the full model. It turns out that the values that maximize the constrained log likelihood are what one should use for the test, and the particular function in which these nuisance parameters are integrated out is termed a profile likelihood.

Consider a parametric model with log likelihood function $\ell(\boldsymbol{\theta})$ whose p -dimensional parameter vector $\boldsymbol{\theta} = (\boldsymbol{\psi}, \boldsymbol{\lambda})$ that can be decomposed into a q -dimensional parameter of interest $\boldsymbol{\psi}$ and a $(p - q)$ -dimensional nuisance vector $\boldsymbol{\lambda}$.

We can consider the profile likelihood ℓ_p , a function of $\boldsymbol{\psi}$ alone, which is obtained by maximizing the likelihood pointwise at each fixed value $\boldsymbol{\psi}_0$ over the nuisance vector $\boldsymbol{\varphi}_{\boldsymbol{\psi}_0}$,

$$\ell_p(\boldsymbol{\psi}) = \max_{\boldsymbol{\varphi}} \ell(\boldsymbol{\psi}, \boldsymbol{\varphi}) = \ell(\boldsymbol{\psi}, \hat{\boldsymbol{\varphi}}_{\boldsymbol{\psi}}).$$

Figure 3.3 shows a fictional log likelihood contour plot with the resulting profile curve (in black), where the log likelihood value is mapped to colors. If one thinks of these contours lines as those of a topographic map, the profile likelihood corresponds in this case to walking along the ridge of both mountains along the $\boldsymbol{\psi}$ direction, with the right panel showing the elevation gain/loss.

The maximum profile likelihood estimator behaves like a regular likelihood for most quantities of interest and we can derive test statistics and confidence intervals in the usual way. One famous example of profile likelihood is the Cox proportional hazard covered in Chapter 7.

Example 3.3 (Box–Cox transformation). Sometimes, the assumption of normality of the error doesn't hold. If the data are strictly positive, one can consider a Box–Cox transformation,

$$y_i(\lambda) = \begin{cases} (y^\lambda - 1)/\lambda, & \lambda \neq 0 \\ \ln(y), & \lambda = 0. \end{cases}$$

If we assume that $\mathbf{y}(\lambda) \sim \text{No}(\mathbf{X}\boldsymbol{\beta}, \sigma^2 \mathbf{I}_n)$, then the likelihood of \mathbf{y} is

$$L(\lambda, \boldsymbol{\beta}, \sigma; \mathbf{y}, \mathbf{X}) = (2\pi\sigma^2)^{-n/2} \exp \left[-\frac{1}{2\sigma^2} \{\mathbf{y}(\lambda) - \mathbf{X}\boldsymbol{\beta}\}^\top \{\mathbf{y}(\lambda) - \mathbf{X}\boldsymbol{\beta}\} \right] J(\lambda, \mathbf{y}),$$

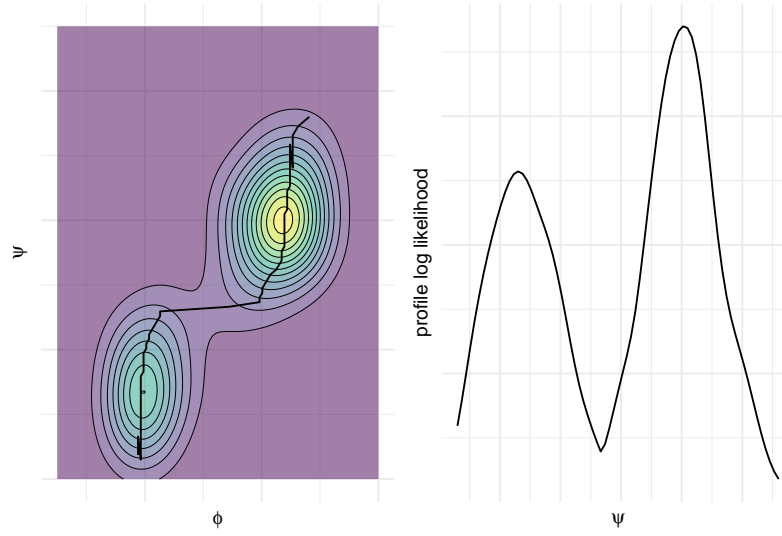


Figure 3.3: Two-dimensional log likelihood surface with a parameter of interest ψ and a nuisance parameter ϕ ; the contour plot shows area of higher likelihood, and the black line is the profile log likelihood, also shown as a function of ψ on the right panel.

where J denotes the Jacobian of the Box–Cox transformation, $\prod_{i=1}^n y_i^{\lambda-1}$. For each given value of λ , the maximum likelihood estimator is that of the usual regression model, with \mathbf{y} replaced by $\mathbf{y}(\lambda)$, namely $\hat{\beta}_\lambda = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}(\lambda)$ and $\hat{\sigma}_\lambda^2 = n^{-1} \{\mathbf{y}(\lambda) - \mathbf{X} \hat{\beta}_\lambda\}^\top \{\mathbf{y}(\lambda) - \mathbf{X} \hat{\beta}_\lambda\}$.

The profile log likelihood is

$$\ell_p(\lambda) = -\frac{n}{2} \ln(2\pi \hat{\sigma}_\lambda^2) - \frac{n}{2} + (\lambda - 1) \sum_{i=1}^n \ln(y_i)$$

The maximum profile likelihood estimator is the value λ minimizes the sum of squared residuals from the linear model with $\mathbf{y}(\lambda)$ as response.

Figure 3.4 shows the profile log likelihood for the linear model with an intercept-only, rescaled to be zero at the maximum. The function shows that a value of approximately 0.37 would provide residuals that are closer to normally distributed. The 95% profile-likelihood based confidence interval is given by the two values of λ , (0.12, 0.69), at which the curve intersects the horizontal grey line drawn at $-\chi_1^2/2$. The Box–Cox is not a panacea and should be reserved to cases where the transformation reduces heteroscedasticity (unequal variance) or creates a linear relation between explanatory and response: theory provides a cogent explanation of the data (e.g., the Cobb–Douglas production function used in economics can be linearized by taking a log-transformation). Rather than an *ad hoc* choice

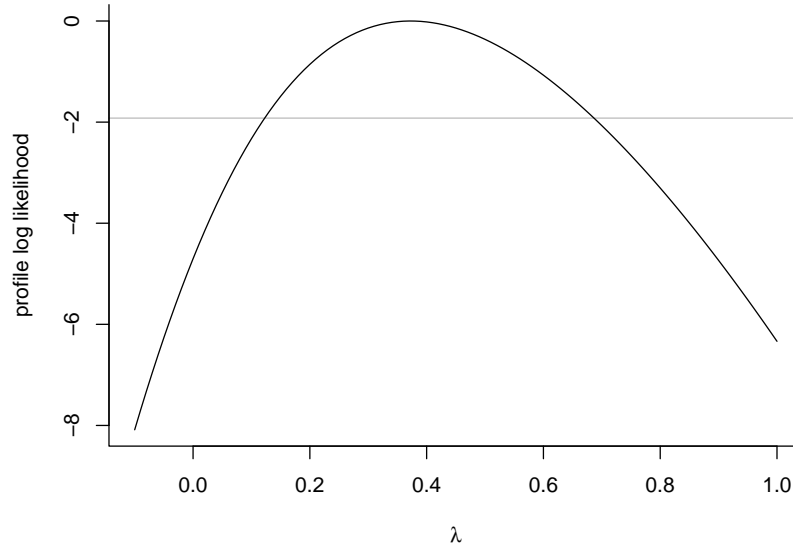


Figure 3.4: Profile log likelihood for the Box–Cox transformation for the waiting time data

of transformation, one could choose a log transformation if the value 0 is included within the 95% confidence interval since this improves interpretability.

We could use the profile likelihood ratio test to obtain confidence intervals for the regression coefficients of the mean model. Consider for simplicity the normal simple linear model, $\mathbf{Y} \sim \text{No}_n(\beta_0 + \beta_1 \mathbf{X}_1, \sigma^2)$. The profile for the slope parameter β_1 can be obtained by maximizing the log likelihood for fixed $\beta_1 = b$, say: to achieve this, note that this amounts to $\mathbf{Y} - b\mathbf{X}_1 \sim \text{No}_n(\beta_0, \sigma^2)$ and so the estimator would correspond to

$$\begin{aligned}\hat{\beta}_0^{(b)} &= \frac{1}{n} \sum_{i=1}^n (Y_i - bX_i) \\ \hat{\sigma}^{2(b)} &= \frac{1}{n} \sum_{i=1}^n (Y_i - bX_i - \hat{\beta}_0^{(b)})^2.\end{aligned}$$

With more than one covariate, we could obtain the value of $\beta_{-j}^{(b)}$ by running least squares and use the residuals to compute the maximum likelihood estimate $\hat{\sigma}^{2(b)}$ of the variance.

3.4 Information criteria

The likelihood can also serve as building block for model comparison: the larger $\ell(\hat{\theta})$, the better the fit. However, the likelihood doesn't account for model complexity in the sense

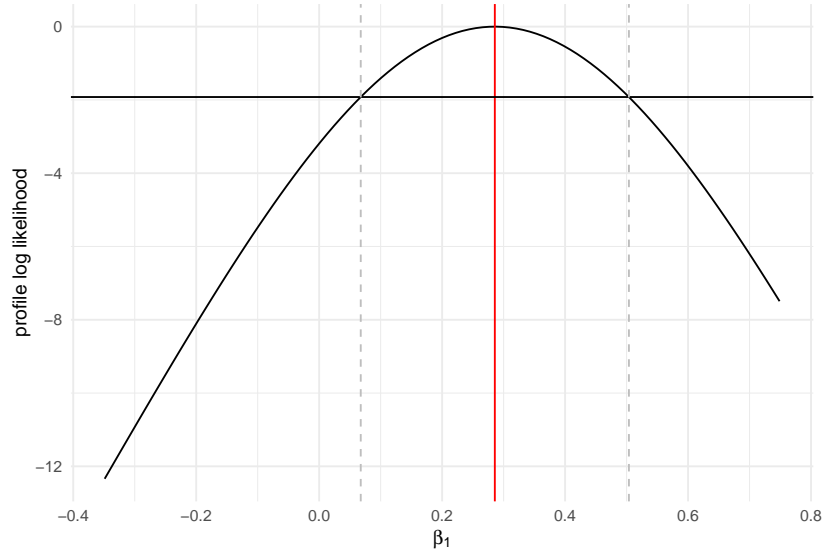


Figure 3.5: Profile log likelihood for β_1 in the simple linear regression model for simulated data. The function has been so its value at the maximum likelihood estimate is zero. The horizontal cutoff marks minus half the 0.95 quantile of the χ_1^2 distribution, with vertical lines indicating the maximum likelihood estimate (red) and 95% confidence interval (dashed gray). Because of the normality assumption of the linear regression model, the sampling distribution is exactly normally distribution, the profile is symmetric and the Wald and profile likelihood ratio based confidence intervals agree.

that more complex models with more parameters lead to higher likelihood. This is not a problem for comparison of nested models using the likelihood ratio test because we look only at relative improvement in fit. There is a danger of **overfitting** if we only consider the likelihood of a model.

AIC and BIC are information criteria measuring how well the model fits the data, while penalizing models with more parameters,

$$\text{AIC} = -2\ell(\hat{\theta}) + 2k$$

$$\text{BIC} = -2\ell(\hat{\theta}) + k \ln(n),$$

where k is the number of parameters in the model. The smaller the value of AIC (or of BIC), the better the model fit.

Note that information criteria do not constitute formal hypothesis tests on the parameters, but they can be used to compare non nested-models, even these estimates are particularly noisy. If we want to compare likelihood from different probability models, we need to

make sure they include normalizing constant. The BIC is more stringent than AIC, as its penalty increases with the sample size, so it selects models with fewer parameters. The BIC is **consistent**, meaning that it will pick the true correct model from an ensemble of models with probability one as $n \rightarrow \infty$. In practice, this is of little interest if one assumes that all models are approximation of reality (it is unlikely that the true model is included in the ones we consider). AIC often selects overly complicated models in large samples, whereas BIC is sometimes too conservative in that it chooses models that are overly simple.

Chapter 4

Generalized linear models

Linear models are only suitable for data that are (approximately) normally distributed. However, there are many settings where we may wish to analyze a response variable which is not necessarily continuous, including when Y is binary, a count variable or is continuous, but non-negative. We will consider in particular likelihood-based inference for binary/proportion and counts data.

Generalized linear models (GLM) combine a model for the conditional mean with a distribution for the response variable and a link function tying predictors and parameters.

This chapter gives an introduction to generalized linear models and focuses in particular on logistic regression and Poisson regression, but only for the case of independent observations. Extensions of generalized linear models for correlated and longitudinal, the so-called *generalized linear mixed models* (GLMM), are covered in MATH80621.

4.1 Basic principles

The starting point is the same as for linear regression: we have a random sample of independent observations (Y, \mathbf{X}) , where Y is the response variable and X_1, \dots, X_p are p explanatory variables or covariates which are assumed fixed (non-random). The goal is to model the mean of the response variable as a function of the explanatory variables.

Let $\mu_i = E(Y_i \mid \mathbf{X}_i)$ denote the conditional expectation of Y_i given covariates and let η_i denote the linear combination of the covariates that will be used to model the response variable,

$$\eta_i = \beta_0 + \beta_1 X_{i1} + \dots + \beta_p X_{ip}.$$

The building blocks of the generalized linear model are

- A random component, consisting of the probability distribution for the outcome Y that is a member of the exponential family (normal, binomial, Poisson, gamma, ...).
- A deterministic component, the linear predictors $\mathbf{X}\beta$, where \mathbf{X} is an $n \times (p + 1)$ matrix with columns $\mathbf{1}_n, X_1, \dots, X_p$ and $\beta \in \mathbb{R}^{p+1}$ are coefficients.
- A monotone function g , called **link function**, that maps the mean of Y_i to the predictor variables, $g(\mu_i) = \eta_i$.

4.2 Theory of generalized linear models

This section borrows from Chapter 4 of

Agresti (2015). * Foundations of Linear and Generalized Linear Models*, Wiley.

In a generalized linear model, the random component arises from an exponential dispersion family: this choice gives a framework, since test statistics and properties can be derived for general classes of distribution.

4.2.1 Exponential-dispersion family of distributions

Consider a probability density or mass function for Y with parameters (θ, ϕ) ,

$$f(y; \theta, \phi) = \exp \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\},$$

where the support, i.e., the set of values taken by Y , doesn't depend on the parameters. Throughout, we will assume **natural parameter** θ is unknown, but the **dispersion parameter** ϕ may be known (exponential family) or unknown (exponential dispersion family).

One particularity of exponential dispersion models is the explicit mean-variance relationship: the first and second derivative of the log likelihood ℓ of a one-sample with respect to the natural parameter θ are

$$\begin{aligned} \frac{\partial \ell(y; \theta, \phi)}{\partial \theta} &= \frac{\partial}{\partial \theta} \left\{ \frac{y\theta - b(\theta)}{a(\phi)} + c(y, \phi) \right\} = \frac{y - b'(\theta)}{a(\phi)} \\ \frac{\partial^2 \ell(y; \theta, \phi)}{\partial \theta^2} &= -\frac{b''(\theta)}{a(\phi)}, \end{aligned}$$

where $b'(\cdot)$ and $b''(\cdot)$ are the first two derivatives of $b(\cdot)$ with respect to θ . Under regularity condition, the Bartlett identities hold and

$$\mathbb{E} \left\{ \frac{\partial \ell(y; \theta, \phi)}{\partial \theta} \right\} = 0, \quad -\mathbb{E} \left\{ \frac{\partial^2 \ell(y; \theta, \phi)}{\partial \theta^2} \right\} = \left[\mathbb{E} \left\{ \frac{\partial \ell(y; \theta, \phi)}{\partial \theta} \right\} \right]^2.$$

These two equality give

$$\begin{aligned} E(Y_i) &= b'(\theta_i) \\ \text{Va}(Y_i) &= b''(\theta_i)a(\phi_i) \end{aligned}$$

Because of the relation between the mean of Y_i , say μ_i , and the natural parameter θ_i , we have $V(\mu) = b''(\theta)$ and there is an explicit relationship between the mean and the variance parameters unless $V(\mu) = 1$.

Example 4.1 (Poisson distribution as exponential family member). The mass function of the Poisson distribution is

$$f(y; \lambda) = \frac{\lambda^y \exp(-\lambda)}{y!} = \exp \{y \ln(\lambda) - \lambda - \ln(y!)\}, \quad y = 0, 1, \dots$$

The natural parameter is $\theta = \ln(\lambda)$, the dispersion parameter $\phi = 1$, and $b(\theta) = \exp(\theta)$. Replacing these expressions in the mean-variance formulas, we get $E(Y) = \exp(\theta) = \mu$ and $\text{Va}(Y) = \exp(\theta) = \mu$, meaning $V(\mu) = \mu$.

Example 4.2 (Binomial distribution as member of the exponential family). We consider the mass function of a binomial distribution $\text{Bin}(m, \pi)$ with the number of trials m known. The parametrization presented in A.3 is not convenient because the mean and the variance both depend on m . We consider thus a different parametrization in which Y represents the fraction of successes, so the mass function takes values in $\{0, 1/m, \dots, 1\}$ and mY denotes the number of successes. The mass function for Y is then

$$\begin{aligned} f(y, \pi) &= \exp \left\{ my \ln \left(\frac{\pi}{1-\pi} \right) + m \ln(1-\pi) + \ln \left[\binom{m}{my} \right] \right\} \\ &= \exp \left\{ \frac{y \ln \left(\frac{\pi}{1-\pi} \right) + \ln(1-\pi)}{1/m} + \ln \left[\binom{m}{my} \right] \right\} \end{aligned}$$

Set

$$\theta = \ln \left(\frac{\pi}{1-\pi} \right)$$

with $b(\theta) = \ln\{1 + \exp(\theta)\}$ and $\phi = m^{-1}$. The expectation and variance are easily derived and

$$\begin{aligned} E(Y) &= \pi = \text{expit}(\theta) = \frac{\exp(\theta)}{1 + \exp(\theta)} = \mu \\ \text{Va}(Y) &= \frac{\pi(1-\pi)}{m} = \frac{\mu(1-\mu)}{m} = \phi V(\mu) \end{aligned}$$

where $V(\mu) = \mu(1-\mu)$.

Example 4.3 (Normal distribution as member of the exponential family). We consider a sample consisting of independent normal observations, $Y_i \sim \text{No}(\mu_i, \sigma^2)$, with

$$\begin{aligned} f(y_i, \mu_i, \sigma^2) &= (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{(y_i - \mu_i)^2}{2\sigma^2} \right\} \\ &= \exp \left[\frac{y_i\mu_i - \mu_i^2/2}{\sigma^2} - \frac{1}{2} \left\{ \frac{y_i^2}{\sigma^2} + \ln(2\pi\sigma^2) \right\} \right], \end{aligned}$$

meaning $\theta = \mu$, $\phi = \sigma^2$ and $a(\phi) = \phi$, $b(\theta) = \theta^2/2$.

4.2.2 Link functions

The link between the mean of Y and the **linear predictor** η is

$$g\{E(Y \mid X_1, \dots, X_p)\} = \eta = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p.$$

The link function thus connects the mean of the random variable Y to the explanatory variables, $g(\mu_i) = \eta_i = \beta_0 + \sum_{j=1}^p \beta_j X_{ij}$ and likewise $\mu_i = g^{-1}(\eta_i)$. If the link function is chosen such that $\theta = \eta$, the link function is termed the canonical link function.

The need for a link function arises from parameter constraints: for example, the mean $\mu = \pi$ of a Bernoulli distribution is the proportion of successes and must lie in the interval $(0, 1)$. Similarly, the mean μ of the Poisson distribution must be positive. For the normal distribution of the ordinary linear regression model, we do not impose constraints on the mean μ_i , so an appropriate link function is the identity, $\mu_i = \eta_i$.

An appropriate choice of link function g sets μ_i equal to a transformation of the linear combination η_i so as to avoid imposing parameter constraints on β . Certain choices of link functions facilitate interpretation or make the likelihood function convenient for optimization.

For the Bernoulli and binomial distributions, an appropriate link function g is the logit function,

$$\begin{aligned} \text{logit}(\mu) &= \ln \left(\frac{\mu}{1 - \mu} \right) = \ln(\mu) - \ln(1 - \mu) = \eta, \\ \text{expit}(\eta) &= \frac{\exp(\eta)}{1 + \exp(\eta)} = \mu. \end{aligned}$$

The inverse link function is the distribution function of the logistic distribution, hence the name. The choice of link function is far from unique: any quantile function of a continuous random variable supported on \mathbb{R} could be considered. For the Poisson distribution, the canonical link function g is the natural logarithm, \ln , with associated inverse link function \exp .

Canonical link functions are natural choices because of their nice statistical properties: choosing the canonical link ensures that $\mathbf{X}^\top \mathbf{y}$ is a minimal sufficient statistic. Other considerations, such as parameter constraints, can be more important in deciding on the choice of g .

Chapter 5

Correlated and longitudinal data

Chapter 6

Linear mixed models

Chapter 7

Survival analysis

Appendix A

Additional topics and prerequisites

A.1 Population and samples

Statistics is the science of uncertainty quantification: of paramount importance is the notion of randomness. Generally, we will seek to estimate characteristics of a population using only a sample (a sub-group of the population of smaller size).

The **population of interest** is a collection of individuals which the study targets. For example, the Labour Force Survey (LFS) is a monthly study conducted by Statistics Canada, who define the target population as “all members of the selected household who are 15 years old and older, whether they work or not.” Asking every Canadian meeting this definition would be costly and the process would be long: the characteristic of interest (employment) is also a snapshot in time and can vary when the person leaves a job, enters the job market or become unemployed.

In general, we therefore consider only **samples** to gather the information we seek to obtain. The purpose of **statistical inference** is to draw conclusions about the population, but using only a share of the latter and accounting for sources of variability. George Gallup made this great analogy between sample and population:

One spoonful can reflect the taste of the whole pot, if the soup is well-stirred

A **sample** is a random sub-group of individuals drawn from the population. Creation of sampling plans is a complex subject and semester-long sampling courses would be required to even scratch the surface of the topic. Even if we won't be collecting data, keep in mind the following information: for a sample to be good, it must be representative of the population under study. Selection bias must be avoided, notably samples of friends or of people sharing opinions.

Because the individuals are selected at **random** to be part of the sample, the measurement of the characteristic of interest will also be random and change from one sample to the next. However, larger samples of the same quality carry more information and our estimator will be more precise. Sample size is not guarantee of quality, as the following example demonstrates.

Example A.1. *The Literary Digest* surveyed 10 millions people by mail to know voting preferences for the 1936 USA Presidential Election. A sizeable share, 2.4 millions answered, giving Alf Landon (57%) over incumbent President Franklin D. Roosevelt (43%). The latter nevertheless won in a landslide election with 62% of votes cast, a 19% forecast error. Biased sampling and differential non-response are mostly responsible for the error: the sampling frame was built using “phone number directories, drivers’ registrations, club memberships, etc.’”, all of which skewed the sample towards rich upper class white people more susceptible to vote for the GOP.

In contrast, Gallup correctly predicted the outcome by polling (only) 50K inhabitants. Read the full story [here](#).

A.2 Random variable

Suppose we wish to describe the behaviour of a stochastic phenomenon. To this effect, one should enumerate the set of possible values taken by the variable of interest and their probability: this is what is encoded in the distribution. We will distinguish between two cases: discrete and continuous variables. Random variables are denoted using capital letters: for example $Y \sim \text{No}(\mu, \sigma^2)$ indicates that Y follows a normal distribution with parameters μ and σ^2 , which represent respectively the expectation and variance of Y .

The (cumulative) distribution function $F(y)$ gives the cumulative probability that an event doesn’t exceed a given numerical value y , $F(y) = \Pr(Y \leq y)$.

If Y is discrete, then it has atoms of non-zero probability and the mass function $f(y) = \Pr(Y = y)$ gives the probability of each outcome y . In the continuous case, no numerical value has non-zero probability and so we consider intervals instead: the density function gives the probability of Y falling in a set B , via $\Pr(Y \in B) = \int_B f(y)dy$. It follows that the distribution function of a continuous random variable is simply $F(y) = \int_{-\infty}^y f(x)dx$.

A.2.1 Moments

One of the first topics covered in introductory statistics is descriptive statistics such as the mean and standard deviation. These are estimators of (centered) moments, which characterise a random variable. In the case of the standard normal distribution, the expectation and variance fully characterize the distribution.

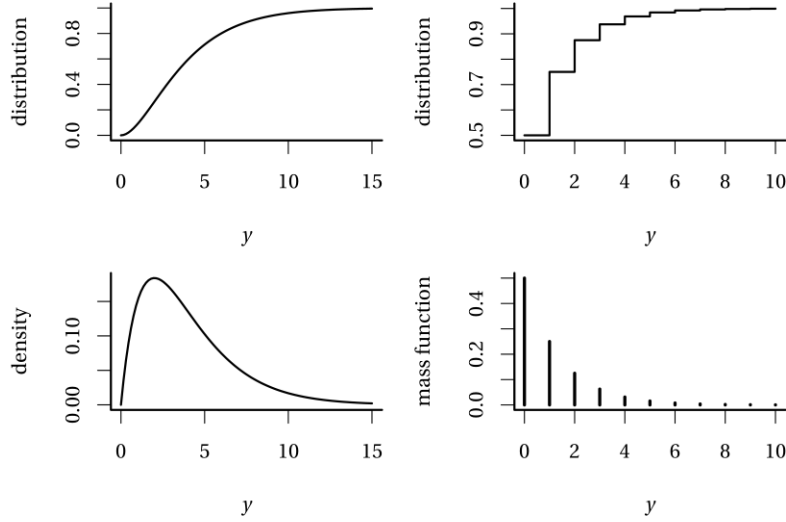


Figure A.1: (Cumulative) distribution functions (top) and density/mass functions (bottom) of continuous (left) and discrete (right) random variables.

Let Y be a random variable with density (or mass) function $f(x)$. This function is non-negative and satisfies $\int_{\mathbb{R}} f(x)dx = 1$: the integral over a set B gives the probability of Y falling inside $B \in \mathbb{R}$.

The expectation (or theoretical mean) of a continuous random variable Y is

$$E(Y) = \int_{\mathbb{R}} x f(x) dx.$$

In the discrete case, we set rather $\mu = E(Y) = \sum_{x \in \mathcal{X}} x \Pr(X = x)$, where \mathcal{X} denotes the support of Y , the set of numerical values at which the probability of Y is non-zero. More generally, we can look at the expectation of a function $g(x)$ for Y , which is nothing but the integral (or sum in the discrete case) of $g(x)$ weighted by the density or mass function of $f(x)$. In the same fashion, provided the integral is finite, the variance is

$$\text{Va}(Y) = E\{Y - E(Y)\}^2 \equiv \int_{\mathbb{R}} (x - \mu)^2 f(x) dx.$$

The **standard deviation** is the square root of the variance and measures the variability of the variable, measured in the same units as Y .

The notion of moments can be extended to higher dimensions. If we assume data are independent, their joint distribution factorizes into the product of the individual mass function/density of each sample point. Otherwise, a multivariate mass function or density function describes the behaviour of the random vector.

Consider an n -vector \mathbf{Y} . In the regression setting, the response \mathbf{Y} would usually comprise repeated measures on an individual, or even observations from a group of individuals.

The expected value (theoretical mean) of the vector \mathbf{Y} is calculated componentwise, i.e.,

$$\mathbb{E}(\mathbf{Y}) = \boldsymbol{\mu} = \left(\mathbb{E}(Y_1) \quad \cdots \quad \mathbb{E}(Y_n) \right)^\top$$

whereas the second moment of \mathbf{Y} is encoded in the $n \times n$ **covariance** matrix

$$\text{Va}(\mathbf{Y}) = \boldsymbol{\Sigma} = \begin{pmatrix} \text{Va}(Y_1) & \text{Co}(Y_1, Y_2) & \cdots & \text{Co}(Y_1, Y_n) \\ \text{Co}(Y_2, Y_1) & \text{Va}(Y_2) & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ \text{Co}(Y_n, Y_1) & \text{Co}(Y_n, Y_2) & \cdots & \text{Va}(Y_n) \end{pmatrix}$$

The i th diagonal element of $\boldsymbol{\Sigma}$, $\sigma_{ii} = \sigma_i^2$, is the variance of Y_i , whereas the off-diagonal entries $\sigma_{ij} = \sigma_{ji}$ ($i \neq j$) are the covariance of pairwise entries, with

$$\text{Co}(Y_i, Y_j) = \int_{\mathbb{R}^2} (y_i - \mu_i)(y_j - \mu_j) f_{Y_i, Y_j}(y_i, y_j) dy_i dy_j.$$

The covariance matrix $\boldsymbol{\Sigma}$ is thus symmetric. It is customary to normalize the pairwise dependence so they do not depend on the component variance. The linear **correlation** between Y_i and Y_j is

$$\rho_{ij} = \text{Cor}(Y_i, Y_j) = \frac{\text{Co}(Y_i, Y_j)}{\sqrt{\text{Va}(Y_i)}\sqrt{\text{Va}(Y_j)}} = \frac{\sigma_{ij}}{\sigma_i \sigma_j}.$$

The correlation matrix of \mathbf{Y} is an $n \times n$ symmetric matrix with ones on the diagonal and the pairwise correlations off the diagonal,

$$\text{Cor}(\mathbf{Y}) = \begin{pmatrix} 1 & \rho_{12} & \rho_{13} & \cdots & \rho_{1n} \\ \rho_{21} & 1 & \rho_{23} & \cdots & \rho_{2n} \\ \rho_{31} & \rho_{32} & 1 & \ddots & \rho_{3n} \\ \vdots & \vdots & \ddots & \ddots & \vdots \\ \rho_{n1} & \rho_{n2} & \rho_{n3} & \cdots & 1 \end{pmatrix}.$$

One of the most important parts of modelling correlated (or longitudinal) data is the need to account for within-group correlations. This basically comes down to modelling a covariance matrix for observations within the same group (or within the same individual in the case of repeated measures), which is the object of Chapter 5.

A.2.2 Unbiasedness and mean square error

An estimator $\hat{\theta}$ for a parameter θ is unbiased if its bias $\text{bias}(\hat{\theta}) = E(\hat{\theta}) - \theta$ is zero. The unbiased estimator of the mean and the variance of Y are

$$\bar{Y}_n = n^{-1} \sum_{i=1}^n Y_i$$

$$S_n = (n-1)^{-1} \sum_{i=1}^n (Y_i - \bar{Y})^2.$$

While unbiasedness is a desirable property, there may be cases where no unbiased estimator exists for a parameter! Often, rather, we seek to balance bias and variance: recall that an estimator is a function of random variables and thus it is itself random: even if it is unbiased, the numerical value obtained will vary from one sample to the next. We often seek an estimator that minimises the mean squared error,

$$\text{MSE}(\hat{\theta}) = E\{(\hat{\theta} - \theta)^2\} = \text{Va}(\hat{\theta}) + \{E(\hat{\theta}) - \theta\}^2.$$

The mean squared error is an objective function consisting of the sum of the squared bias and the variance.

An alternative to this criterion is to optimize a function such as the likelihood of the sample: the resulting estimator is termed maximum likelihood estimator. These estimator are asymptotically efficient, in the sense that they have the lowest mean squared error of all estimators for large samples. Other properties of maximum likelihood estimators make them attractive default choice for estimation.

The likelihood of a sample is the joint density of the n observations, which requires a distribution to be considered. Many such distributions describe simple physical phenomena and can be described using a few parameters: we only cover the most frequently encountered.

A.2.3 Discrete distributions

Example A.2 (Bernoulli distribution). We consider a binary event such as coin toss (heads/tails). In general, the two events are associated with success/failure. By convention, failures are denoted by zeros and successes by ones, the probability of success being π so $\Pr(Y = 1) = \pi$ and $\Pr(Y = 0) = 1 - \pi$ (complementary event). The mass function of the Bernoulli distribution is thus

$$\Pr(Y = y) = \pi^y (1 - \pi)^{1-y}, \quad y = 0, 1.$$

A rapid calculation shows that $E(Y) = \pi$ and $\text{Va}(Y) = \pi(1 - \pi)$. Many research questions have binary responses, for example:

- did a potential client respond favourably to a promotional offer?
- is the client satisfied with service provided post-purchase?
- will a company go bankrupt in the next three years?
- did a study participant successfully complete a task?

Example A.3 (Binomial distribution). If the data give the sum of independent Bernoulli events, the number of successes Y out of m trials is binomial, denoted $\text{Bin}(m, \pi)$; the mass function of the binomial distribution is

$$\Pr(Y = y) = \binom{m}{y} \pi^y (1 - \pi)^{1-y}, \quad y = 0, 1.$$

The likelihood of a sample from a binomial distribution is (up to a normalizing constant that doesn't depend on π) the same as that of m independent Bernoulli trials. The expectation of the binomial random variable is $E(Y) = m\pi$ and its variance $\text{Va}(Y) = m\pi(1 - \pi)$.

As examples, we could consider the number of successful candidates out of m who passed their driving license test or the number of customers out of m total which spent more than 10\$ in a store.

More generally, we can also consider count variables whose realizations are integer-valued, for examples the number of

- insurance claims made by a policyholder over a year,
- purchases made by a client over a month on a website,
- tasks completed by a study participant in a given time frame.

Example A.4 (Geometric distribution). The geometric distribution is a model describing the number of Bernoulli trials with probability of success π required to obtain a first success. The mass function of $Y \sim \text{Geo}(\pi)$ is

$$\Pr(Y = y) = \pi(1 - \pi)^{y-1}, \quad y = 1, 2, \dots$$

For example, we could model the numbers of visits for a house on sale before the first offer is made using a geometric distribution.

Example A.5 (Poisson distribution). If the probability of success π of a Bernoulli event is small in the sense that $m\pi \rightarrow \lambda$ when the number of trials m increases, then the number of success follows approximately a Poisson distribution with mass function

$$\Pr(Y = y) = \frac{\exp(-\lambda)\lambda^y}{\Gamma(y + 1)}, \quad y = 0, 1, 2, \dots$$

where $\Gamma(\cdot)$ denotes the gamma function. The parameter λ of the Poisson distribution is both the expectation and the variance of the distribution, meaning $E(Y) = \text{Va}(Y) = \lambda$.

Example A.6 (Negative binomial distribution). The negative binomial distribution arises as a natural generalization of the geometric distribution if we consider the number of Bernoulli trials with probability of success π until we obtain m success. Let Y denote the number of failures: the order of success and failure doesn't matter, but for the latest trial which is a success. The mass function is thus

$$\Pr(Y = y) = \binom{m-1+y}{y} \pi^m (1-\pi)^y.$$

The negative binomial distribution also appears as the unconditional distribution of a two-stage hierarchical gamma-Poisson model, in which the mean of the Poisson distribution is random and follows a gamma distribution. In notation, this is $Y \mid \Lambda = \lambda \sim \text{Po}(\lambda)$ and Λ follows a gamma distribution with shape r and scale θ , whose density is

$$f(x) = \theta^{-r} x^{r-1} \exp(-x/\theta) / \Gamma(r).$$

The unconditional number of success is then negative binomial.

In the context of generalized linear models, we will employ yet another parametrisation of the distribution, with the mass function

$$\Pr(Y = y) = \frac{\Gamma(y+r)}{\Gamma(y+1)\Gamma(r)} \left(\frac{r}{r+\mu} \right)^r \left(\frac{\mu}{r+\mu} \right)^y, y = 0, 1, \dots, \mu, r > 0,$$

where Γ is the gamma function and the parameter $r > 0$ is not anymore integer valued. The expectation and variance of Y are $E(Y) = \mu$ et $\text{Va}(Y) = \mu + k\mu^2$, where $k = 1/r$. The variance of the negative binomial distribution is thus higher than its expectation, which justifies the use of the negative binomial distribution for modelling overdispersion.

A.2.4 Continuous distributions

We will encounter many continuous distributions that arise as (asymptotic) null distribution of test statistics.

Example A.7 (Normal distribution). The normal distribution is ubiquitous in statistics because of the central limit theorem. A random variable Y follows a normal distribution if its density function is

$$f(x) = (2\pi\sigma^2)^{-1/2} \exp \left\{ -\frac{(x-\mu)^2}{2\sigma^2} \right\}, \quad x \in \mathbb{R}.$$

The parameters μ and $\sigma > 0$ that fully characterize the distribution of the normal distribution and they correspond to the expectation and standard deviation. The normal

distribution is a location-scale distribution, so $(Y - \mu)/\sigma \sim \text{No}(0, 1)$. The distribution function of the standard normal distribution, Φ , is not available in closed-form.

We will also encounter the multivariate normal distribution; for a n dimensional vector $\mathbf{Y} \sim \text{No}_p(\boldsymbol{\mu}, \boldsymbol{\Sigma})$, the density is

$$f(\mathbf{x}) = (2\pi)^{-p/2} |\boldsymbol{\Sigma}|^{-1/2} \exp \left\{ -\frac{1}{2} (\mathbf{x} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1} (\mathbf{x} - \boldsymbol{\mu}) \right\}$$

The mean vector $\boldsymbol{\mu}$ is the vector of expectation of individual observations, whereas $\boldsymbol{\Sigma}$ is the covariance matrix of \mathbf{Y} . A unique property of the multivariate normal distribution is the link between independence and the covariance matrix: if Y_i and Y_j are independent, the (i, j) off-diagonal entry of $\boldsymbol{\Sigma}$ is zero.

Example A.8 (Chi-square distribution). The chi-square distribution arises as null distribution of likelihood-based test. If $\mathbf{Y} \sim \text{No}_p(\mathbf{0}_p, \mathbf{I}_p)$, i.e., all components are independent and centered $Y_i \sim \text{No}(0, 1)$, then $\sum_{i=1}^p Y_i^2$ follows a chi-square distribution with p degrees of freedom, denote χ_p^2 . The square of a standard normal variate likewise follows a χ_1^2 distribution.

If we consider a sample of n normally distributed observations, the sample variance $(n - 1)S^2/\sigma^2 \sim \chi_{n-1}^2$.

Example A.9 (Student- t distribution). If $X \sim \text{No}(0, 1)$ independent of $Y \sim \chi_\nu^2$, then

$$T = \frac{X}{\sqrt{Y/\nu}}$$

follows a Student- t distribution with ν degrees of freedom, denoted St_ν . The density of T is

$$f(y; \nu) = \frac{\Gamma\left(\frac{\nu+1}{2}\right)}{\Gamma\left(\frac{\nu}{2}\right) \sqrt{\nu\pi}} \left(1 + \frac{y^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

the distribution has polynomial tails, is symmetric around 0 and unimodal. As $\nu \rightarrow \infty$, the Student distribution converges to a normal distribution. It has heavier tails than the normal distribution and only the first $\nu - 1$ moments of the distribution exist, so a Student distribution with $\nu = 2$ degrees of freedom has infinite variance.

For normally distributed data, the centered sample mean divided by the sample variance, $(\bar{Y} - \mu)/S^2$ follows a Student- t distribution with $n - 1$ degrees of freedom, which explains the terminology t -tests.

A.2.5 Quantiles-quantiles plots

Models are (at best) an approximation of the true data generating mechanism and we will want to ensure that our assumptions are reasonable and the quality of the fit decent. Quantile-quantile plots are graphical goodness-of-fit diagnostics that are based on the following principle: if Y is a continuous random variable with distribution function F , then the mapping $F(Y) \sim U(0, 1)$ yields uniform variables. Similarly, the quantile transform applied to a uniform variable provides a mean to simulating samples from F , viz. $F^{-1}(U)$. Consider then a random sample of size n from the uniform distribution ordered from smallest to largest, with $U_{(1)} \leq \dots \leq U_{(n)}$. One can show these ranks have marginally a Beta distribution, $U_{(k)} \sim \text{Beta}(k, n + 1 - k)$ with expectation $k/(n + 1)$.

In practice, we don't know F and, even if we did, one would need to estimate the parameters. We consider some estimator \hat{F} for the model and apply the inverse transform to an approximate uniform sample $\{i/(n + 1)\}_{i=1}^n$. The quantile-quantile plot shows the data as a function of the (first moment) of the transformed order statistics:

- on the x -axis, the theoretical quantiles $\hat{F}^{-1}\{\text{rank}(y_i)/(n + 1)\}$
- on the y -axis, the empirical quantiles y_i

If the model is adequate, the ordered values should follow a straight line with unit slope passing through the origin. Whether points fall on a 45 degree line is difficult to judge by eye and so it is advisable to ease the interpretation to subtract the slope: the detrended plot is easier to interpret and was proposed by Tukey (but beware of the scale of the y -axis!). Figure A.2 shows two representations of the same data using simulated samples from a standard normal distribution.

Even if we knew the true distribution of the data, the sample variability makes it very difficult to spot if deviations from the model are abnormal or compatible with the model. A simple point estimate with no uncertainty measure can lead to wrong conclusions. As such, we add approximate pointwise or simultaneous confidence intervals. The simplest way to do this is by simulation (using a parametric bootstrap), by repeating the following steps B times:

1. simulate a (bootstrap) sample $\{Y_i^{(b)}\}(i = 1, \dots, n)$ from \hat{F}
2. re-estimate the parameters of F to obtain $\hat{F}_{(b)}$
3. calculate and save the plotting positions $\hat{F}_{(b)}^{-1}\{i/(n + 1)\}$.

The result of this operation is an $n \times B$ matrix of simulated data. We obtain a symmetric $(1 - \alpha)$ confidence interval by keeping the empirical quantile of order $\alpha/2$ and $1 - \alpha/2$ from each row. The number B should be larger than 999, say, and be chosen so that B/α is an integer.

For the pointwise interval, each order statistic from the sample is a statistic and so the

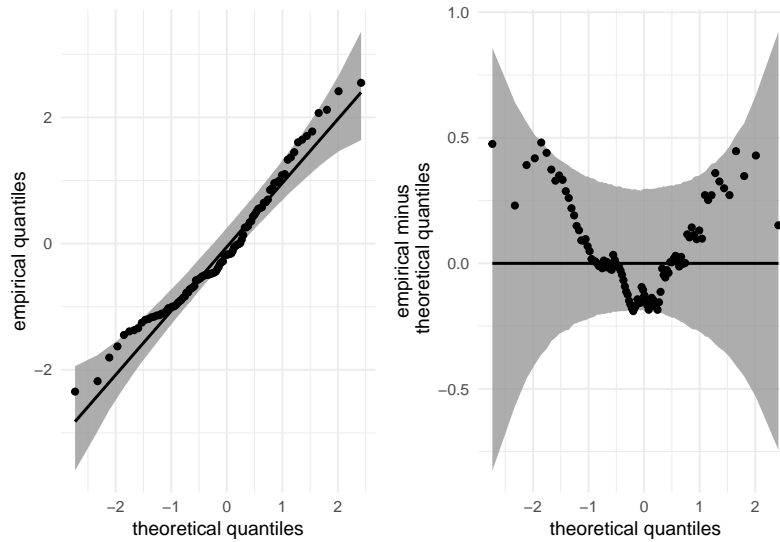


Figure A.2: Normal quantile-quantile plot (left) and detrended version (Tukey's) of the same plot (right).

probability of any single one falling outside the confidence interval is approximately α . However, order statistics are not independent (they are ordered), so it's common to see neighboring points falling outside of their respective intervals. [It is also possible to use the bootstrap samples to derive an (approximate) simultaneous confidence intervals, in which we expected values to fall $100(1 - \alpha)\%$ of the time inside the bands in repeated samples; see Section 4.4.3 of these course notes. The intervals shown in Figure A.2 are pointwise and derived (magically) using a simple function. The uniform order statistics have larger variability as we move away from 0.5, but the uncertainty in the quantile-quantile plot largely depends on F .

Interpretation of quantile-quantile plots requires some experience: this post by *Glen_b* on StackOverflow nicely summarizes what can be detected (or not) from them.

A.3 Laws of large numbers

An estimator for a parameter θ is **consistent** if the value obtained as the sample size increases (to infinity) converges to the true value of θ . Mathematically speaking, this translates into convergence in probability, meaning $\hat{\theta} \xrightarrow{\text{Pr}} \theta$. In common language, we say that the probability that $\hat{\theta}$ and θ differ becomes negligible as n gets large.

Consistency is the *a minima* requirement for an estimator: when we collect more informa-

tion, we should approach the truth. The law of large number states that the sample mean of n (independent) observations with common mean μ , say \bar{Y}_n , converges to μ , denoted $\bar{Y}_n \rightarrow \mu$. Roughly speaking, our approximation becomes less variable and asymptotically unbiased as the sample size (and thus the quantity of information available for the parameter) increases. The law of large number is featured in Monte Carlo experiments: we can approximate the expectation of some (complicated) function $g(x)$ by simulating repeatedly independent draws from Y and calculating the sample mean $n^{-1} \sum_{i=1}^n g(Y_i)$.

If the law of large number tells us what happens in the limit (we get a single numerical value), the result doesn't contain information about the rate of convergence and the uncertainty at finite levels.

A.4 Central Limit Theorem

The central limit theorem is perhaps the flagship result of probability theory: for a random sample of size n with (independent) random variables whose expectation is μ and variance σ^2 , then the sample mean converges to μ , but

- the estimator \bar{Y} is centered around μ ,
- the standard error is σ/\sqrt{n} ; the rate of convergence is thus \sqrt{n} . For a sample of size 100, the standard error of the sample mean will be 10 times smaller than that of the underlying random variable.
- the sample mean, once properly scaled, follows approximately a normal distribution

Mathematically, the central limit theorem states $\sqrt{n}(\bar{Y} - \mu) \xrightarrow{d} \text{No}(0, \sigma^2)$. If n is large (a rule of thumb is $n > 30$, but this depends on the underlying distribution of Y), then $\bar{Y} \sim \text{No}(\mu, \sigma^2/n)$.

How do we make sense of this result? Let us consider the mean travel time of high speed Spanish trains (AVE) between Madrid and Barcelona that are operated by Renfe.

Our exploratory data analysis showed previously that the duration is the one advertised on the ticket: there are only 15 unique travel time. Based on 9603 observations, we estimate the mean travel time to be 170 minutes and 41 seconds. Figure A.3 shows the empirical distribution of the data.

Consider now samples of size $n = 10$, drawn repeatedly from the population: in the first sample, the sample mean is 170.9 minutes, whereas we get an estimate of 164.5 minutes in our second, 172.3 minutes in the third, etc.

We draw $B = 1000$ different samples, each of size $n = 5$, from two millions records, and calculate the sample mean in each of them. The top right panel of A.4 shows the result for $n = 5$, but also for $n = 20$ (bottom left). The last graph of Figure A.4 shows the impact of

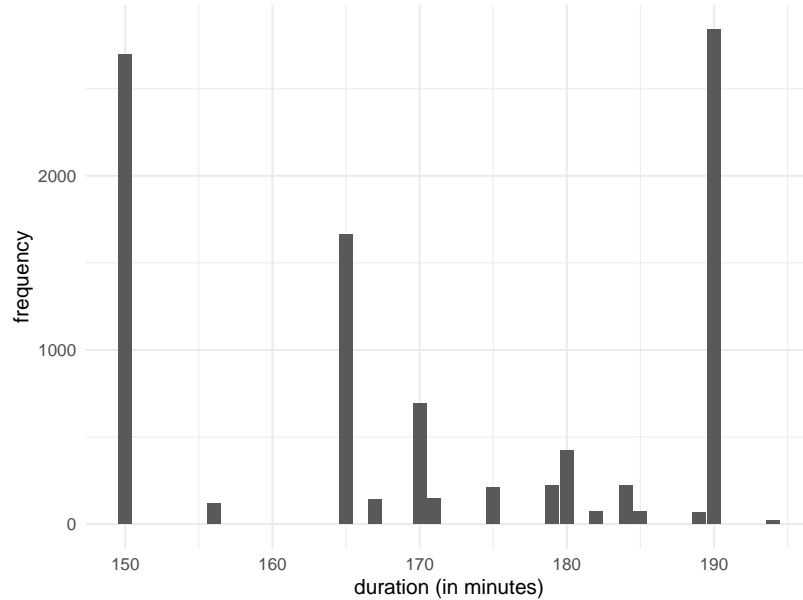


Figure A.3: Empirical distribution of travel times of high speed trains.

the increase in sample size: whereas the normal approximation is okay-ish for $n = 5$, it is indistinguishable from the normal approximation for $n = 20$. As n increases and the sample size gets bigger, the quality of the approximation improves and the curve becomes more concentrated around the true mean. Even if the distribution of the travel time is discrete, the mean is approximately normal.

We considered a single distribution in the example, but you could play with other distributions and vary the sample size to see when the central limit theorem kicks in using this applet.

The central limit theorem underlies why scaled test statistics which have sample mean zero and sample variance 1 have a standard null distribution in large sample: this is what guarantees the validity of our inference!

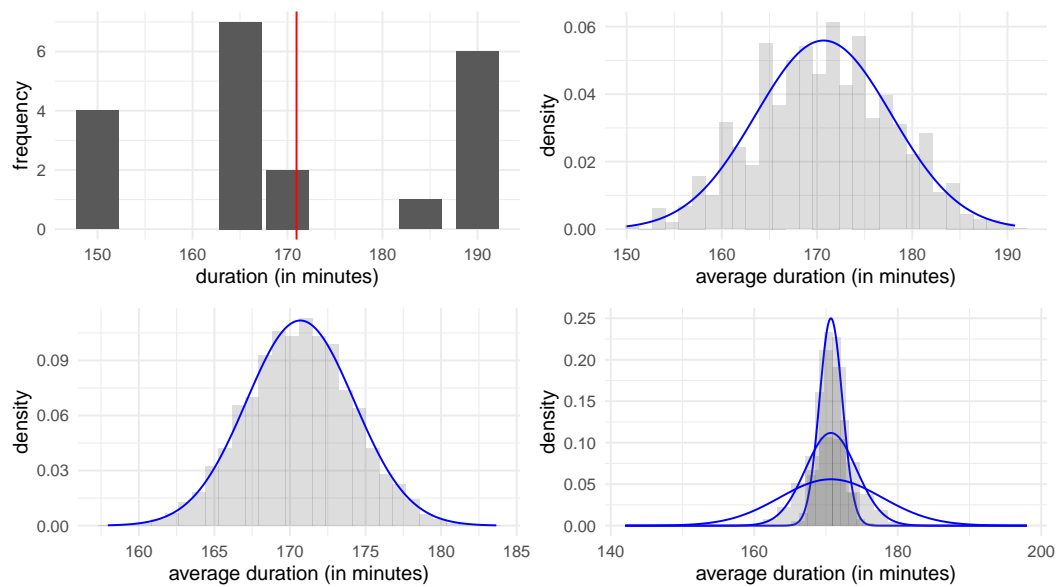


Figure A.4: Graphical representation of the central limit theorem. The upper left panel shows a sample of 20 observations with its sample mean (vertical red). The three other panels show the histograms of the sample mean from repeated samples of size 5 (top right), 20 (bottom left) and 20, 50 and 100 overlaid, with the density approximation provided by the central limit theorem.

Appendix B

Mathematical derivations

This section regroups optional derivations which are provided for the sake of completeness.

B.1 Derivation of the ordinary least squares estimator

Consider the optimization problem

$$\hat{\beta} = \min_{\beta \in \mathbb{R}^{p+1}} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta).$$

We can compute the derivative of the right hand side with respect to β , set it to zero and solve for $\hat{\beta}$,

$$\begin{aligned} \mathbf{0}_n &= \frac{\partial}{\partial \beta} (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \\ &= \frac{\partial (\mathbf{y} - \mathbf{X}\beta)}{\partial \beta} \frac{\partial (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta)}{\partial (\mathbf{y} - \mathbf{X}\beta)} \\ &= \mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) \end{aligned}$$

using the chain rule. Distributing the terms leads to the so-called *normal equation*

$$\mathbf{X}^\top \mathbf{X} \beta = \mathbf{X}^\top \mathbf{y}.$$

If the $n \times p$ matrix \mathbf{X} is full-rank, the quadratic form $\mathbf{X}^\top \mathbf{X}$ is invertible and we obtain the solution to the least square problems provided in Equation (2.3).

B.2 Derivation of the coefficient of determination

Because of the orthogonal decomposition $\mathbf{y} = \hat{\mathbf{y}} + \mathbf{e}$ and provided that the design matrix includes an intercept of $\mathbf{1}_n \in \mathcal{S}(\mathbf{X})$, then $\bar{e} = 0$ and the average of the response and of the fitted values is the same. Since $n^{-1} \sum_{i=1}^n \hat{y}_i = n^{-1} \sum_{i=1}^n (y_i - e_i) = \bar{y}$,

$$\begin{aligned}
 \widehat{\text{Cor}}(\hat{\mathbf{y}}, \mathbf{y}) &= \frac{(\mathbf{y} - \bar{y}\mathbf{1}_n)^\top (\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n)}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\| \|\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n\|} \\
 &= \frac{(\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n)^\top (\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n) + \mathbf{e}^\top (\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n)}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\| \|\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n\|} \\
 &= \frac{\|\hat{\mathbf{y}} - \bar{y}\mathbf{1}_n\|}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\|} \\
 &= \frac{\|\mathbf{y} - \bar{y}\mathbf{1}_n\| - \|\mathbf{e}\|}{\|\mathbf{y} - \bar{y}\mathbf{1}_n\|} \\
 &= \sqrt{\frac{\text{SS}_c - \text{SS}_e}{\text{SS}_c}} = R.
 \end{aligned}$$

This justifies the claim of Section 2.5 that the squared correlation between the fitted values and the response is equal to R^2 .

Appendix C

R

R is an object-oriented interpreted language. It differs from usual programming languages in that it is designed for interactive analyses.

You can find several introductions to **R** online. Have a look at the **R** manuals or better at contributed manuals. A nice official reference is *An introduction to R*. You may wish to look up the following chapters of the **R** language definition (Evaluation of expressions and part of the *Objects* chapter). Another good (small reference) is the cheatsheet *Getting started in R*.

C.1 Basics of R

Help

Help can be accessed via `help` or simply `?`, e.g., `help("Normal")`. See **R** page about help files.

Basic commands

Basic **R** commands are fairly intuitive, especially if you want to use **R** as a calculator. Elementary functions such as `sum`, `min`, `max`, `sqrt`, `log`, `exp`, etc., are self-explanatory.

Some (unconventional) features of the language:

- **R** is case sensitive.
- Use `<-` for assignments to a variable, and `=` for matching arguments inside functions
- Indexing in **R** starts at 1, **not** 0.
- Most functions in **R** are vectorized and loops are typically inefficient.

- Integers are obtained by appending L to the number, so 2L is an integer and 2 a double (numerical).

Besides integers and doubles, the common types are

- logical (TRUE and FALSE);
- null pointers (NULL), which can be assigned to arguments;
- missing values, namely NA or NaN. These can also be obtained a result of invalid mathematical operations such as $\log(-2)$.

Beware! In **R**, invalid calls will often returns *something* rather than an error. It is therefore good practice to check that the output is sensical.

Linear algebra in **R**

R is an object oriented language, and the basic elements in **R** are (column) vector. Below is a glossary with some useful commands for performing basic manipulation of vectors and matrix operations:

- `c` concatenates elements to form a vector
- `cbind` (`rbind`) binds column (row) vectors
- `matrix` and `vector` are constructors
- `diag` creates a diagonal matrix (by default with ones)
- `t` is the function for transpose
- `rep` creates a vector of duplicates, `seq` a sequence. For integers i, j with $i < j$, `i:j` generates the sequence $i, i + 1, \dots, j - 1, j$.

Subsetting is fairly intuitive and general; you can use vectors, logical statements. For example, if `x` is a vector, then

- `x[2]` returns the second element
- `x[-2]` returns all but the second element
- `x[1:5]` returns the first five elements
- `x[(length(x) - 5):length(x)]` returns the last five elements
- `x[c(1, 2, 4)]` returns the first, second and fourth element
- `x[x > 3]` return any element greater than 3. Possibly an empty vector of length zero!
- `x[x < -2 | x > 2]` multiple logical conditions.
- `which(x == max(x))` index of elements satisfying a logical condition.

For a matrix `x`, subsetting now involves dimensions: `[1, 2]` returns the element in the first row, second column. `x[, 2]` will return all of the rows, but only the second column. For lists, you can use `[[` for subsetting by index or the `$` sign by names.

Packages

The great strength of **R** comes from its contributed libraries (called packages), which contain functions and datasets provided by third parties. Some of these (base, stats, graphics, etc.) are loaded by default whenever you open a session.

To install a package from CRAN, use `install.packages("package")`, replacing `package` by the package name. Once installed, packages can be loaded using `library(package)`; all the functions in `package` will be available in the environment.



There are drawbacks to loading packages: if an object with the same name from another package is already present in your environment, it will be hidden. Use the double-colon operator `::` to access a single object from an installed package (`package::object`).

Datasets

- datasets are typically stored inside a `data.frame`, a matrix-like object whose columns contain the variables and the rows the observation vectors.
- The columns can be of different types (integer, double, logical, character), but all the column vectors must be of the same length.
- Variable names can be displayed by using `names(faithful)`.
- Individual columns can be accessed using the column name using the `$` operator. For example, `faithful$eruptions` will return the first column of the `faithful` dataset.
- To load a dataset from an (installed) **R** package, use the command `data` with the name of the package as an argument (must be a string). The package `datasets` is loaded by default whenever you open **R**, so these are always in the search path.

The following functions can be useful to get a quick glimpse of the data:

- `summary` provides descriptive statistics for the variable.
- `str` provides the first few elements with each variable, along with the dimension
- `head` (`tail`) prints the first (last) n lines of the object to the console (default is $n = 6$).

We start by loading a dataset of the Old Faithful Geyser of Yellowstone National park and looking at its entries.

```
# Load Old faithful dataset
data(faithful, package = "datasets")
# Query the database for documentation
?faithful
```

```
# look at first entries
head(faithful)
```

```
##   eruptions waiting
## 1      3.60      79
## 2      1.80      54
## 3      3.33      74
## 4      2.28      62
## 5      4.53      85
## 6      2.88      55
```

```
str(faithful)
```

```
## 'data.frame': 272 obs. of 2 variables:
## $ eruptions: num 3.6 1.8 3.33 2.28 4.53 ...
## $ waiting : num 79 54 74 62 85 55 88 85 51 85 ...
```

```
# What kind of object is faithful?
class(faithful)
```

```
## [1] "data.frame"
```

Other common classes of objects:

- `matrix`: an object with attributes `dim`, `ncol` and `nrow` in addition to `length`, which gives the total number of elements.
- `array`: a higher dimensional extension of `matrix` with arguments `dim` and `dimnames`.
- `list`: an unstructured class whose elements are accessed using double indexing `[[]]` and elements are typically accessed using `$` symbol with names. To delete an element from a list, assign `NULL` to it.
- `data.frame` is a special type of list where all the elements are vectors of potentially different type, but of the same length.

Base graphics

The `faithful` dataset consists of two variables: the regressand `waiting` and the regressor `eruptions`. One could postulate that the waiting time between eruptions will be smaller if the eruption time is small, since pressure needs to build up for the eruption to happen. We can look at the data to see if there is a linear relationship between the variables.

An image is worth a thousand words and in statistics, visualization is crucial. Scatterplots are produced using the function `plot`. You can control the graphic console options using `par` — see `?plot` and `?par` for a description of the basic and advanced options available.

Once `plot` has been called, you can add additional observations as points (lines) to the graph using `point` (`lines`) in place of `plot`. If you want to add a line (horizontal, vertical, or with known intercept and slope), use the function `abline`.

Other functions worth mentioning for simple graphics:

- `boxplot` creates a box-and-whiskers plot
- `hist` creates an histogram, either on frequency or probability scale (option `freq = FALSE`). `breaks` control the number of bins. `rug` adds lines below the graph indicating the value of the observations.
- `pairs` creates a matrix of scatterplots, akin to `plot` for data frame objects.



There are two options for basic graphics: the base graphics package and the package `ggplot2`. The latter is a more recent proposal that builds on a modular approach and is more easily customizable — I suggest you stick to either and `ggplot2` is a good option if you don't know **R** already, as the learning curve will be about the same. Even if the display from `ggplot2` is nicer, this is no excuse for not making proper graphics. Always label the axis and include measurement units!

C.2 Linear models in **R** using the *lm* function

The function `lm` is the workhorse for fitting linear models in **R**. It takes as input a formula: suppose you have a data frame containing columns `x` (a regressor) and `y` (the regressand); you can then call `lm(y ~ x)` to fit the linear model $y = \beta_0 + \beta_1 x + \varepsilon$. The explanatory variable `y` is on the left hand side, while the right hand side should contain the predictors, separated by a `+` sign if there are more than one. If you provide the data frame name using `data`, then the shorthand `y ~ .` fits all the columns of the data frame (but `y`) as regressors.

If you include categorical variables, make sure they are transformed to factors. Normally, strings are cast to factor (unless you specify `stringsAsFactors = FALSE`) upon import of the data, but the danger here lies with variables that are encoded using integers (sex, revenue class, level of education, marital status, etc.) It is okay if we keep binary variables as is if they are encoded using 0/1, but it is often better to cast them to factor to get more meaningful labels given the lack of obvious ordering.

By default, the baseline level for a factor is based on the alphabetical order, while **SAS** uses the first value it encounters. Once the variables are cast to factor, `summary` will print the counts for each respective categories; these could be likewise be obtained using `table`.

To fit higher order polynomials, use `poly`. For transformations, use the `I` function to tell **R** to interpret the input “as is”. Thus, `lm(y~x+I(x^2))`, would fit a linear model with design matrix $(1_n, \mathbf{x}^\top, \mathbf{x}^2)^\top$. A constant is automatically included in the regression, but can be removed by writing `-1` or `+0` on the right hand side of the formula (but don’t do that!). The `lm` function output will display ordinary least squares estimates along with standard errors, t values for the Wald test of the hypothesis $H_0 : \beta_i = 0$ and the associated P -values. Other statistics and information about the sample size, the degrees of freedom, etc., are given at the bottom of the table.

Many methods allow you to extract specific objects from `lm` objects. For example, the functions `coef`, `resid`, `fitted`, `model.matrix` will return the coefficients $\hat{\beta}$, the ordinary residuals e , the fitted values \hat{y} and the design matrix \mathbf{X} , respectively.

```
data(college, package = "hecstatmod") #load data
class(college$rank) #check that rank is a factor
```

```
## [1] "factor"
```

```
#if not, use the following "<-" to assign, "$" to access the column of a data.frame/list
# college$rank <- factor(college$rank, labels = c("assistant","associate","full"))
linmod <- lm(salary ~ sex + rank + service + field, data = college)
coef(linmod) #coefficients
```

```
##      (Intercept)      sexwoman      rankassociate      rankfull
##      86.5963      -4.7712      14.5604      49.1596
##      service fieldtheoretical
##      -0.0888      -13.4734
```

```
summary(linmod) #summary table
```

```
##
## Call:
## lm(formula = salary ~ sex + rank + service + field, data = college)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -64.20 -14.26  -1.53   10.57   99.16
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
```

```
## (Intercept)      86.5963      2.9603      29.25 < 2e-16 ***
## sexwoman        -4.7712      3.8780      -1.23  0.21931
## rankassociate    14.5604      4.0983       3.55  0.00043 ***
## rankfull         49.1596      3.8345      12.82 < 2e-16 ***
## service         -0.0888      0.1116      -0.80  0.42696
## fieldtheoretical -13.4734      2.3155      -5.82  1.2e-08 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 22.7 on 391 degrees of freedom
## Multiple R-squared:  0.448, Adjusted R-squared:  0.441
## F-statistic: 63.4 on 5 and 391 DF,  p-value: <2e-16
```

```
confint(linmod) #confidence intervals for model parameters
```

```
##              2.5 % 97.5 %
## (Intercept)  80.776 92.416
## sexwoman     -12.396  2.853
## rankassociate  6.503 22.618
## rankfull      41.621 56.698
## service       -0.308  0.131
## fieldtheoretical -18.026 -8.921
```

```
yhat <- fitted(linmod) #fitted values
e <- resid(linmod) #ordinary residuals
```


Bibliography

- Brockwell, P. & Davis, R. (2016). *Introduction to Time Series and Forecasting*. Springer Texts in Statistics. Springer.
- Davison, A. C. (2008). *Statistical Models*. Cambridge Series in Statistical and Probabilistic Mathematics. Cambridge: Cambridge University Press.
- Fox, J. & Monette, G. (1992). Generalized collinearity diagnostics. *Journal of the American Statistical Association*, 87(417), 178–183.
- Satterthwaite, F. E. (1946). An approximate distribution of estimates of variance components. *Biometrics Bulletin*, 2(6), 110–114.
- Welch, B. L. (1947). The generalization of “Student’s” problem when several population variances are involved. *Biometrika*, 34(1–2), 28–35.