

# Practical Statistics for Experimental Biologists in RStudio

Caroline Walter

August 2020

## Contents

<b>1</b>	<b>Background</b>	<b>2</b>
<b>2</b>	<b>Contingency Tables</b>	<b>2</b>
<b>3</b>	<b>Sensitivity and Specificity</b>	<b>4</b>
<b>4</b>	<b>Positive Predictive Value (PPV)</b>	<b>4</b>
<b>5</b>	<b>Negative Predictive Value (NPV)</b>	<b>5</b>
<b>6</b>	<b>Prevalence</b>	<b>6</b>
<b>7</b>	<b>Frequentist vs. Bayesian Statistics</b>	<b>7</b>
<b>8</b>	<b>Fisher's Exact Test</b>	<b>7</b>
<b>9</b>	<b>Bernoulli Distribution</b>	<b>8</b>
<b>10</b>	<b>Binomial Tests</b>	<b>8</b>
<b>11</b>	<b>Chi Square Tests</b>	<b>9</b>
<b>12</b>	<b>Confidence Intervals and P-values</b>	<b>10</b>
<b>13</b>	<b>Gaussian/Normal Distributions</b>	<b>10</b>
<b>14</b>	<b>Standard Error of the Mean (SEM)</b>	<b>10</b>
<b>15</b>	<b>Z-Tests</b>	<b>11</b>
<b>16</b>	<b>T-Tests</b>	<b>12</b>
<b>17</b>	<b>QQ Plots</b>	<b>13</b>
<b>18</b>	<b>The Correlation Coefficient (r)</b>	<b>15</b>
<b>19</b>	<b>The Coefficient of Determination (r squared)</b>	<b>15</b>
<b>20</b>	<b>Power Analysis</b>	<b>16</b>
<b>21</b>	<b>ANOVA tests (Analysis of Variance)</b>	<b>18</b>
<b>22</b>	<b>Nonlinear Regression</b>	<b>20</b>

# 1 Background

This course will be reviewing some practical statistics work in the R computing language useful in experimental biology.

Below are the necessary libraries for the statistical work we will be doing.

```
#Loading libraries
library(caret)
library(tidyverse)
library(dplyr)
library(RDocumentation)
library(psych)
library(car)
library(pwr)
```

# 2 Contingency Tables

- create an example data set

```
example_data <- data.frame(x1 = sample(letters[1:5], 20, replace=TRUE),
                           x2 = sample(LETTERS[1:5], 20, replace = TRUE))
```

- have a look at the data set

```
print.data.frame(example_data)
```

```
##      x1 x2
## 1    e  C
## 2    b  C
## 3    a  E
## 4    b  B
## 5    d  B
## 6    c  B
## 7    a  E
## 8    e  C
## 9    d  E
## 10   d  B
## 11   d  A
## 12   d  C
## 13   d  E
## 14   d  C
## 15   c  C
## 16   e  E
## 17   c  E
## 18   d  D
## 19   b  B
## 20   b  A
```

- create a table

```
my_table_0 <- table(example_data$x1, example_data$x2)
print.table(my_table_0)
```

```
##
##      A B C D E
## a 0 0 0 0 2
## b 1 2 1 0 0
## c 0 1 1 0 1
## d 1 2 2 1 2
## e 0 0 2 0 1
```

- if we want to have row and column totals

```
my_table_01<- addmargins(my_table_0)
print.table(my_table_01)
```

```
##
##      A  B  C  D  E Sum
## a      0  0  0  0  2   2
## b      1  2  1  0  0   4
## c      0  1  1  0  1   3
## d      1  2  2  1  2   8
## e      0  0  2  0  1   3
## Sum    2  5  6  1  6  20
```

- convert it to a dataframe

```
my_table_1 <- as.data.frame.matrix(my_table_0)
```

- have a look at the table

```
print.data.frame(my_table_1)
```

```
##      A B C D E
## a 0 0 0 0 2
## b 1 2 1 0 0
## c 0 1 1 0 1
## d 1 2 2 1 2
## e 0 0 2 0 1
```

- to have a table of proportions based on rows, and convert it to a dataframe

```
my_table_2 <- prop.table(my_table_0, margin = 1) %>%
  as.data.frame.matrix()
```

- have a look at the table

```
print.data.frame(my_table_2, digits = 2)
```

```
##      A      B      C      D      E
## a 0.00 0.00 0.00 0.00 1.00
## b 0.25 0.50 0.25 0.00 0.00
## c 0.00 0.33 0.33 0.00 0.33
## d 0.12 0.25 0.25 0.12 0.25
## e 0.00 0.00 0.67 0.00 0.33
```

- to have a table of proportions based on columns

```
my_table_3 <- prop.table(my_table_0, margin=2)%>%
  as.data.frame.matrix()
```

- have a look at the table

```
print.data.frame(my_table_2, digits = 2)
```

```
##      A    B    C    D    E
## a 0.00 0.00 0.00 0.00 1.00
## b 0.25 0.50 0.25 0.00 0.00
## c 0.00 0.33 0.33 0.00 0.33
## d 0.12 0.25 0.25 0.12 0.25
## e 0.00 0.00 0.67 0.00 0.33
```

### 3 Sensitivity and Specificity

For the AIP test: AIP (Acute Intermediate Porphyria) is a rare metabolic disorder.

- Sensitivity: (AIP example) the probability of testing positive given that the subject has the disease

$$\text{Sensitivity} = p(\text{test}^+ | \text{disease}^+) = 0.82$$

$$p(\text{test}^+ | \text{disease}^+) = \frac{TP}{TP+FN} = \frac{82}{82+18} = 0.82$$

- Specificity: (AIP example) the probability of a negative test given that the subject does not have the disease

$$\text{Specificity} = p(\text{test}^- | \text{disease}^-) = 0.963$$

- here is a helpful youtube video: <https://www.youtube.com/watch?v=9f5XgjWpzi0>

```
devtools::install_github("datacamp/RDocumentation")
```

```
data01 <- factor(c("A", "B", "B", "B"))
data02 <- factor(c("A", "B", "B", "B"))

ref01 <- factor(c("B", "B", "B", "B"))
ref02 <- factor(c("B", "A", "B", "B"))
```

### 4 Positive Predictive Value (PPV)

Positive Predictive Value is the probability that subjects with a positive screening test truly have the disease

- TP: True Positive
- FP: False Positive
- TN: True Negative
- FN: False Negative

$$p(\text{disease}^+ | \text{test}^+) = \frac{TP}{TP+FP} = \frac{82}{82+36,996} = 0.0022$$

- Even if you test positive, the probability of you having AIP is still very low.
- PPV is often far less than sensitivity in screening tests for rare diseases.

```
table(data01, ref01)
```

```
##      ref01
## data01 B
##      A 1
##      B 3
```

```

sensitivity(data01, ref01)

## [1] 0.75
posPredValue(data01, ref01)

## [1] NA
table(data02, ref02)

##          ref02
## data02 A B
##      A 0 1
##      B 1 2
sensitivity(data02, ref02)

## [1] 0
posPredValue(data02, ref02)

## [1] 0
data03 <- factor(c("A", "B", "B", "B"))
data04 <- factor(c("B", "B", "B", "B"))

ref03 <- factor(c("B", "B", "B", "B"), levels = c("A", "B"))
ref04 <- factor(c("B", "A", "B", "B"))

```

## 5 Negative Predictive Value (NPV)

Negative predictive value is the probability that subjects with a negative screening test truly do not have the disease.

$$p(\text{disease}^- | \text{test}^-)$$

```

table(data03, ref03)

##          ref03
## data03 A B
##      A 0 1
##      B 0 3
specificity(data03, ref03) #0.75

## [1] 0.75
negPredValue(data03, ref03)

## [1] NA
table(data04, ref04)

##          ref04
## data04 A B
##      B 1 3
specificity(data04, ref04) #1

## [1] 1

```

```

negPredValue(data04, ref04) #NaN

## [1] NaN

if(!isTRUE(all.equal(sensitivity(data01, ref01), .75))) stop("error in sensitivity test 1")
if(!isTRUE(all.equal(sensitivity(data02, ref02), 0))) stop("error in sensitivity test 2")
ref03 <- factor(c("B", "B", "B", "B"))
if(!is.na(sensitivity(data02, ref03, "A"))) stop("error in sensitivity test3")

options(show.error.messages = FALSE)
test1 <- try(sensitivity(data02, as.character(ref03)))
if(grep("Error", test1) != 1)
  stop("error in sensitivity calculation - allowed non-factors")
options(show.error.messages = TRUE)

ref03 <- factor(c("B", "B", "B", "B"), levels = c("A", "B"))

if(!isTRUE(all.equal(specificity(data03, ref03), .75))) stop("error in specificity test 1")

if(!isTRUE(all.equal(specificity(data04, ref04), 1.00))) stop("error in specificity test 2")

if(!is.na(specificity(data01, ref01, "A"))) stop("error in specificity test3")

options(show.error.messages = FALSE)
test1 <- try(specificity(data04, as.character(ref03)))
if(grep("Error", test1) != 1)
  stop("error in specificity calculation - allowed non-factors")
options(show.error.messages = TRUE)

```

## 6 Prevalence

Prevalence is the fraction of individuals in a population who have a disease.

- If a subject's sibling has AIP, there is a 50% chance that they do too.  $prevalence = p(disease^+) = 0.50$
- Usage: `prevalence(model, type = c("pop", "bnp", "wnp"), i=NULL, ...)`

```

set.seed(0934)
Data.All.df.2008 <- data.frame(FSA = sample(c("N8N", "N8R", "B3L", "P1H"), 50, T),
                              Lyme = sample(0:1, 50, T),
                              stringsAsFactors = F)

```

- First 10 observations:

```
head(Data.All.df.2008)
```

```

##   FSA Lyme
## 1 N8N    1
## 2 P1H    1
## 3 N8N    0
## 4 P1H    0
## 5 N8N    1
## 6 N8N    1

```

- Prevalence can be calculated as the number of positive diagnoses divided by the total number of observations, i.e. `sum(Lyme)/n()`.

```
Data.All.df.2008 %>%
  group_by(FSA) %>%
  summarise(Prevalence = sum(Lyme)/n())
```

```
## # A tibble: 4 x 2
##   FSA   Prevalence
##   <chr>      <dbl>
## 1 B3L        0.778
## 2 N8N        0.571
## 3 N8R        0.583
## 4 P1H        0.467
```

## 7 Frequentist vs. Bayesian Statistics

- Frequentist Statistics (aka classical statistics) focusses on likelihood. It avoids calculations involving prior odds, and therefore yields results that are prone to misinterpretation due to the base rate fallacy. Frequentist statistics is used heavily in biological research. It can still be useful and informative if you know exactly what to watch out for.

$p(\text{data}|\text{hypothesis})$

- Base Rate Fallacy: If presented with related base rate information (generic, general information) and
- Iron Law of Frequentist Statistics: Never compute the probability of a hypothesis.
- Bayesian Statistics explicitly accounts for prior odds. It focuses on computing posterior probabilities and requires prior information that is often hard to quantify. It is central to the modern machine learning and more advanced areas of quantitative biology. Experimental researchers in biology tend not to use Bayesian statistics.

$p(\text{hypothesis}|\text{data})$

## 8 Fisher's Exact Test

Fisher's exact test is a statistical significance test used in the analysis of contingency tables and in place of chi square test in 2 by 2 tables, especially in cases of small samples. The test is useful for categorical data that result from classifying objects in two different ways; it is used to examine the significance of the association (contingency) between the two kinds of classification.

Example: A British woman claimed to be able to distinguish whether milk or tea was added to the cup first. To test, she was given 8 cups of tea, in four of which milk was added first. The null hypothesis is that there is no association between the true order of pouring and the woman's guess, the alternative that there is a positive association (that the odds ratio is greater than 1).

Syntax:

```
fisher.test(x, y = NULL, workspace = 200000, hybrid = FALSE,
  hybridPars = c(expect = 5, percent = 80, Emin = 1),
  control = list(), or = 1, alternative = "two.sided",
  conf.int = TRUE, conf.level = 0.95,
  simulate.p.value = FALSE, B = 2000)
```

```
TeaTasting <-
  matrix(c(3,1,1,3),
        nrow=2,
        dimnames = list(Guess = c("Milk", "Tea"),
                        Truth = c("Milk", "Tea")))
fisher.test(TeaTasting, alternative = "greater")

##
## Fisher's Exact Test for Count Data
##
## data: TeaTasting
## p-value = 0.2429
## alternative hypothesis: true odds ratio is greater than 1
## 95 percent confidence interval:
##  0.3135693      Inf
## sample estimates:
## odds ratio
##  6.408309
```

## 9 Bernoulli Distribution

Bernoulli Distribution: a discrete distribution having two possible outcomes labeled by  $n=0$  and  $n=1$  in which  $n=1$  (“success”) occurs with probability  $p$  and  $n=0$  (“failure”) occurs with probability  $q = 1-p$ , where  $0 < p < 1$ . Describes probabilities for a binary variable. (Biased coin example); when the probabilities sum up to 100%.

Example: A Biased Coin Biased coins are modeled using a Bernoulli distribution, which describes probabilities for a binary variable.

## 10 Binomial Tests

A binomial test compares the number of successes observed in a given number of trials with a hypothesized probability of success. The test has the null hypothesis that the real probability of success is equal to some value denoted  $p$ , and the alternative hypothesis that is not equal to  $p$ . The test can also be performed with a one-sided alternative hypothesis that the real probability of success is either greater than  $p$  or that it is less than  $p$ .

A one-tailed test with a significance level of 0.05 will be used. You roll the die 300 times and throw a total of 60 sixes. We cannot reject the null hypothesis that the probability of rolling a six is  $1/6$ . This means that there is no evidence to prove that the die is not fair.

- Syntax: `binom.test(nsuccesses, ntrials, p)`

```
binom.test(60, 300, 1/6, alternative = "greater")

##
## Exact binomial test
##
## data: 60 and 300
## number of successes = 60, number of trials = 300, p-value = 0.07299
## alternative hypothesis: true probability of success is greater than 0.1666667
## 95 percent confidence interval:
##  0.1626847 1.0000000
## sample estimates:
```



```
## probability of success
## 0.2
```

## 11 Chi Square Tests

Chi Square tests are used to determine if two categorical variables have a significant correlation between them. The two variables are selected from the same population. (male/female, red/green, yes/no). The chi-square statistic is commonly used for testing relationships between categorical variables.

$$\chi^2 = \sum \frac{(\text{observed} - \text{expected})^2}{\text{expected}}$$

- Syntax: `chisq.test(data)`
- importing data, have a look at the table

```
data_frame<- read.csv("https://goo.gl/j6lRXD")
table(data_frame$treatment, data_frame$improvement)
```

```
##
##           improved not-improved
## not-treated      26           29
## treated         35           15
```

- Example A: Here, we get a chi-squared value of 5.5569. Since we get a p-value less than the significance level of 0.05, we reject the null hypothesis and conclude that the two variables are in fact dependent.

```
chisq.test(data_frame$treatment, data_frame$improvement, correct=FALSE)
```

```
##
## Pearson's Chi-squared test
##
## data: data_frame$treatment and data_frame$improvement
## X-squared = 5.5569, df = 1, p-value = 0.01841
```

- Example B: Here, we get a high chi-squared value and a p-value of less than the 0.05 significance level. Therefore, we can reject the null hypothesis and conclude that carb and cyl have a significant relationship.

```
data("mtcars")
table(mtcars$carb, mtcars$cyl)
```

```
##
##      4 6 8
## 1 5 2 0
## 2 6 0 4
## 3 0 0 3
## 4 0 4 6
## 6 0 1 0
## 8 0 0 1
```

```
chisq.test(mtcars$carb, mtcars$cyl)
```

```
## Warning in chisq.test(mtcars$carb, mtcars$cyl): Chi-squared approximation may be
## incorrect
```

```
##
## Pearson's Chi-squared test
##
```

```
## data:  mtcars$carb and mtcars$cyl
## X-squared = 24.389, df = 10, p-value = 0.006632
```

## 12 Confidence Intervals and P-values

- P-values quantify the probability of data being as or more extreme than the data in hand, were the null hypothesis true.
- Confidence Intervals are more informative than p-values. We can reject a null hypothesis if it lies outside of the confidence interval.

## 13 Gaussian/Normal Distributions

Gaussian distribution (“the normal distribution”) is ubiquitous in statistics. This distribution is based on the mean  $\mu$  and standard deviation  $\sigma$  of the data.

$x \sim Normal(\mu, \sigma^2)$

- The Central Limit Theorem states that the population of all possible samples of size  $n$  from a populations= with mean  $\mu$  and variance ( $\sigma$  squared) approaches a normal distribution when  $n$  (sample size) approaches infinity. The central limit theorem makes the normal distribution extremely relevant.

The parameters of a statistical model that has been fit to a large dataset will have lingering uncertainty, but this uncertainty will very often be approximately normally distributed. This is why statisticians so often assume that experimental measurements follow normal distributions.

$\theta$  = model parameter inferred from data  $\theta \sim Normal(\mu, \sigma^2)$

The goal of statistical inference is to determine the mean and standard deviation of this distribution  $\mu$ : best estimate of  $\theta$ , denoted  $\theta$   $\sigma$ : lingering uncertainty in  $\theta$ : affects confidence interval

- Example: Here, this normal/gaussian distribution results shows that the percentage of students scoring an 84% or higher in the college entrance exam is 21.5%

```
pnorm(84, mean = 72, sd = 15.2, lower.tail = FALSE)
```

```
## [1] 0.2149176
```

## 14 Standard Error of the Mean (SEM)

The standard error of the mean (SEM) is a statistical term that measures the accuracy with which a sample distribution represents a population by using the standard deviation. In statistics, a sample mean deviates from the actual mean of a population - this deviation is the standard error of the mean (SEM)

$$SEM = \sqrt{q(1 - q)/N}$$

```
Input =(
Stream      Fish
Mill_Creek_1      76
Mill_Creek_2     102
North_Branch_Rock_Creek_1  12
North_Branch_Rock_Creek_2  39
Rock_Creek_1      55
Rock_Creek_2     93
```

```

Rock_Creek_3      98
Rock_Creek_4      53
Turkey_Branch     102
")

Data = read.table(textConnection(Input),header=TRUE)

```

- calculate standard error manually

```

sd(Data$Fish, na.rm=TRUE) /
sqrt(length(Data$Fish[!is.na(Data$Fish)]))

```

```
## [1] 10.69527
```

- install psych package in console and use describe function from package for standard error. This function also works on dataframes.

```

describe(Data$Fish,
          type=2)

```

```

##      vars n mean      sd median trimmed  mad min max range  skew kurtosis   se
## X1      1 9   70 32.09      76      70 34.1  12 102   90 -0.65    -0.69 10.7

```

## 15 Z-Tests

The z-test function is based on the standard normal distribution and creates confidence intervals and tests hypotheses for both one and two sample problems.

$$z = \frac{x - \mu}{\sigma}$$

An R function called `z.test()` would be great for doing the kind of testing in which you use z-scores in the hypothesis test. It could look like this:

- Syntax: `z.test(x, y=NULL, alternative = "two-sided", mu =  $\theta$ , sigma.x=NULL, sigma.y=NULL, conf.level=0.95)`
- $\mu$  = a single number representing the value of the mean or difference in means specified by the null hypothesis
- $\sigma_x$  = a single number representing the population standard deviation for x
- $\sigma_y$  = a single number representing the population standard deviation for y

One problem: that function does not exist in base R. Although you can find one in other packages, it's easy enough to create one and learn a bit about R programming in the process. The function would work like this:

Example: `* IQ.data <- c(100,101,104, 109, 125, 116, 105, 108, 110) * z.test(IQ.data, 100, 15) * z=1.733`  
 Results: `* one-tailed probability = 0.042 * two-tailed probability = 0.084`

Begin by creating the function name and its arguments: `* z.test=function(x,  $\mu$ , popvar){` \* The first argument is the vector of data, second is the population mean, third is the population variance. Left curly bracket signifies that the remainder of the code is what happens inside the function. Next, create a vector that will hold the one-tailed probability of the z-score you will calculate: `* one.tail.p <- NULL` Then you calculate the z-score and round it to three decimal places `* z.score<- round((mean(x)-mu)/popvar/sqrt(length(x))), 3)`

## 16 T-Tests

A t-test is a statistical test which is used to compare the mean of two groups of samples. It is therefore to evaluate whether the means of the two sets of data are statistically significantly different from each other.

The `t.test()` function produces a variety of t-tests. Unlike most statistical packages, the default assume unequal variance and applies the Welch degrees of freedom modification.

A one-sample t-test is used to compare the mean of a population with a theoretical value. An unpaired two sample t-test is used to compare the mean of two independent samples. A paired t-test is used to compare the means between two related groups of samples. The formula of a t-test depends on the mean and standard deviation of the data being compared.

- One-Sample t-test: based on the theoretical mean( $\mu$ ), the set of values with size n, the mean m, and the standard deviation  $\sigma$ . The degrees of freedom is found with (df=n-1).

$$t = \frac{m - \mu}{\frac{\sigma}{\sqrt{n}}}$$

- Independent two-sample t-test Let A and B represent the two groups to compare. Let mA and mB represent the means of groups A and B, respectively. Let nA and nB represent the sizes of group A and B respectively.

$$t = \frac{mA - mB}{\sqrt{\frac{S^2}{nA} + \frac{S^2}{nB}}}$$

$S^2$  is an estimator of the common variance of the two samples.

- Welch's T-Tests (or unequal variances t-test), is a two-sample location test which is used to test the hypothesis that the two populations have unequal means.
- independent two-sample t-test, where y is a numeric and x is a binary factor

```
x=rnorm(10)
y=rnorm(10)
t.test(y,x) #Welch's t-test
```

```
##
##  Welch Two Sample t-test
##
## data:  y and x
## t = 0.055054, df = 12.403, p-value = 0.957
## alternative hypothesis: true difference in means is not equal to 0
## 95 percent confidence interval:
##  -0.9643282  1.0145094
## sample estimates:
##  mean of x      mean of y
##  0.018022788 -0.007067778
```

- one sample t-test, with a null hypothesis ( $H_0$ )  $\mu = 3$

```
t.test(y,mu=3)
```

```
##
##  One Sample t-test
##
## data:  y
## t = -7.1566, df = 9, p-value = 5.328e-05
## alternative hypothesis: true mean is not equal to 3
## 95 percent confidence interval:
##  -0.9245609  0.9606065
```

```
## sample estimates:
## mean of x
## 0.01802279
```

- Mann-Whitney-Wilcoxon Test Two data samples are independent if they come from distinct populations and the samples do not affect each other. Using the Mann-Whitney-Wilcoxon Test, we can decide whether the population distribution are identical without assuming them to follow the normal distribution.

Example: At the 0.05 significance level, we conclude that the gas mileage data of manual and automatic transmissions in mtcars are nonidentical populations.

```
mtcars$mpg
```

```
## [1] 21.0 21.0 22.8 21.4 18.7 18.1 14.3 24.4 22.8 19.2 17.8 16.4 17.3 15.2 10.4
## [16] 10.4 14.7 32.4 30.4 33.9 21.5 15.5 15.2 13.3 19.2 27.3 26.0 30.4 15.8 19.7
## [31] 15.0 21.4
```

```
wilcox.test(mpg ~ am, data=mtcars)
```

```
## Warning in wilcox.test.default(x = c(21.4, 18.7, 18.1, 14.3, 24.4, 22.8, :
## cannot compute exact p-value with ties

##
## Wilcoxon rank sum test with continuity correction
##
## data: mpg by am
## W = 42, p-value = 0.001871
## alternative hypothesis: true location shift is not equal to 0
```

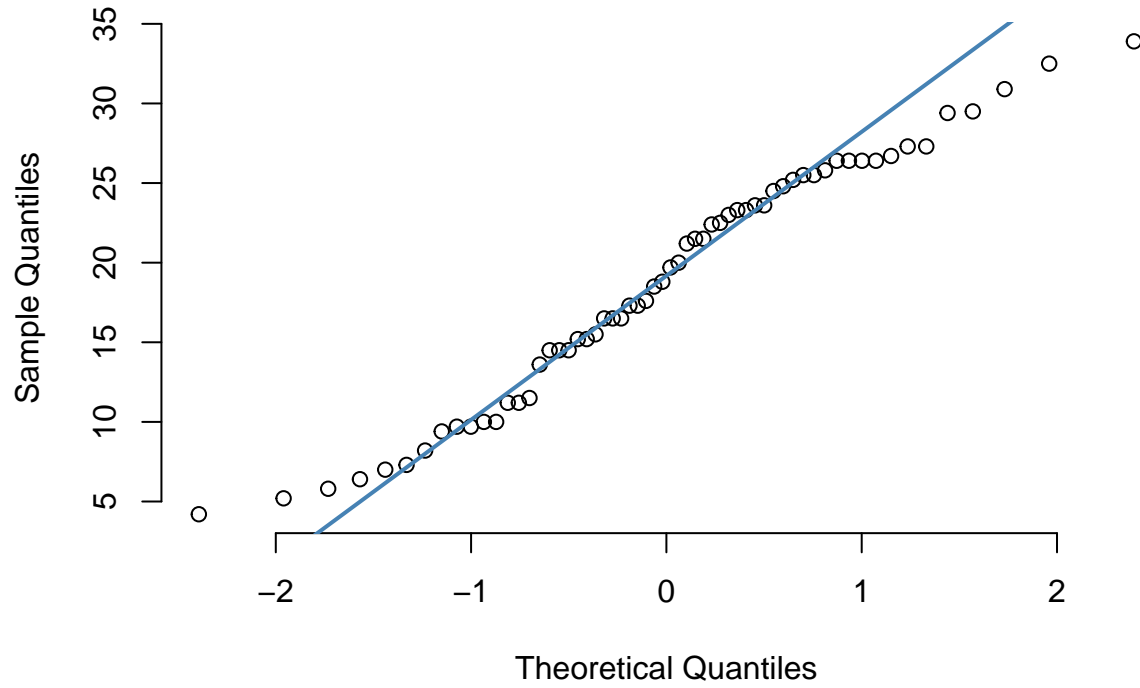
## 17 QQ Plots

QQ Plots are used to visually test whether data follows an expected distribution. They are used to verify that data used in a t-test is actually normally distributed. QQ plots are not useful on small datasets. As an example, we will use data from the `car` package, so install and load `car` (`install.packages("car")` and `library(car)`). Then:

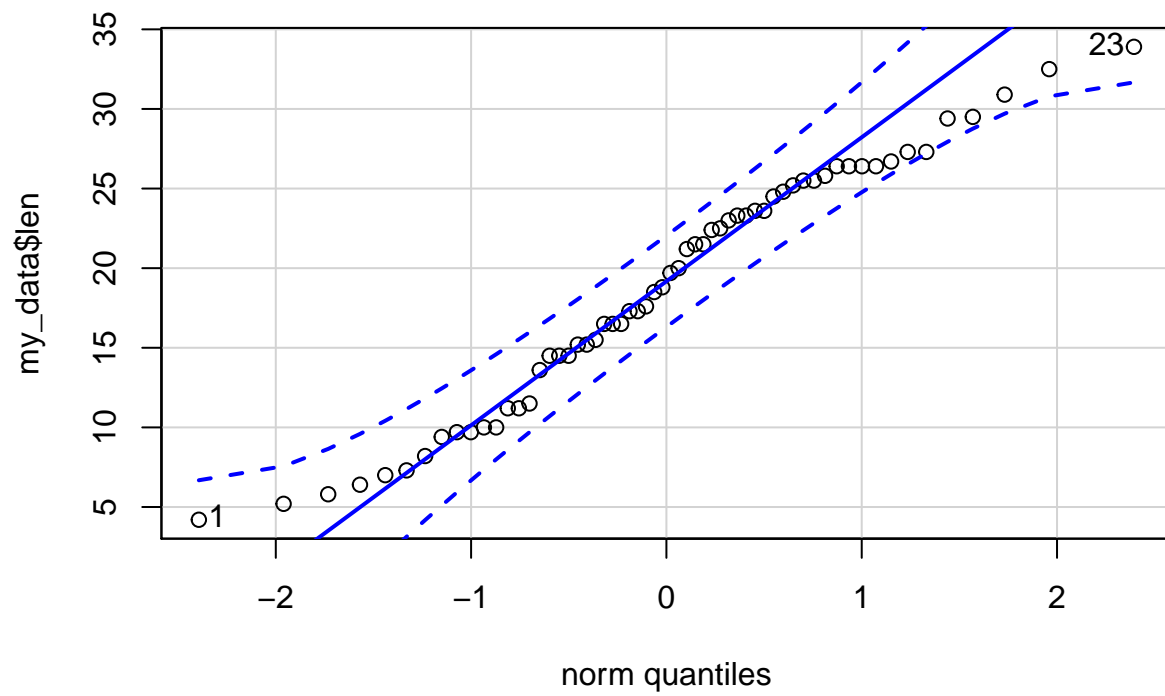
```
my_data <- ToothGrowth
```

```
qqnorm(my_data$len, pch = 1, frame = FALSE)
qqline(my_data$len, col = "steelblue", lwd=2)
```

# Normal Q-Q Plot



```
qqPlot(my_data$len)
```



```
## [1] 23 1
```

## 18 The Correlation Coefficient (r)

The main result of a correlation is called the correlation coefficient (“r”). It ranges from -1 to 1. The closer r is to -1 or 1, the more closely the two variables are (inversely) related. If r is close to 0, it means there is no relationship between the variables. \* When  $r=0$ , the two variables are independent. \* When  $r=+1/-1$ , the two variables share a deterministic linear relationship.

## 19 The Coefficient of Determination (r squared)

R-squared is always between 0 and 1. It is commonly interpreted as the fraction of variance in y explained by x (or the other way around). It is a goodness-of-fit measure for linear regression models. This statistics indicates a percentage of the variance in the dependent variable that the independent variables explain collectively. T-squared measures the strength of the relationship between your model and the dependent variable on a convenient 0-100% scale.

After fitting a linear regression model to the data (by using a least-squares regression line), we need to determine how well the model fits the data.

Assessing Goodness-of-Fit in a Regression Model \* Linear Regression identifies the equation that produces the smallest difference between all of the observed values and their fitted values. Linear regression finds the smallest sum of squared residuals that is possible for the dataset (least-squares regression line).

R-Squared is the percentage of the dependent variable variation that a linear model explains.

$$R^2 = \frac{\text{variance explained by model}}{\text{total variance}}$$

- 0% or an r-squared close to 0.000 represents a model that does not explain any of the variation in the response (dependent) variable around its mean. The mean of the dependent variable predicts the dependent variable as well as the regression model.
- 100% or an r-squared of 1.00 represents a model that explains all of the variation in the response variable around its mean.
- Example: We apply the `lm` function to a formula that describes the variable `eruptions` by the variable `waiting`, and save the linear regression model in a new variable `eruption.lm`.

```
eruption.lm = lm(eruptions ~ waiting, data=faithful)
summary(eruption.lm)$r.squared
```

```
## [1] 0.8114608
```

```
summary(eruption.lm)
```

```
##
## Call:
## lm(formula = eruptions ~ waiting, data = faithful)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.29917 -0.37689  0.03508  0.34909  1.19329
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept) -1.874016   0.160143  -11.70  <2e-16 ***
## waiting      0.075628   0.002219   34.09  <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

```
##
## Residual standard error: 0.4965 on 270 degrees of freedom
## Multiple R-squared:  0.8115, Adjusted R-squared:  0.8108
## F-statistic: 1162 on 1 and 270 DF,  p-value: < 2.2e-16
```

The coefficient of determination (r squared) of the simple linear regression model for the data set below is 0.81146. This means with an r-squared of 0.81146, then approximately 81.15% of the observed variation can be explained by the model's inputs.

## 20 Power Analysis

The power of a statistical test is the probability that the test will reject a false null hypothesis. Power analysis allows us to determine the sample size required to detect an effect of a given size with a given degree of confidence. It also allows us to determine the probability of detecting an effect of a given size with a given level of confidence, under sample size constraints.

- Example: The function tells us we should flip the coin 22.55126 times, which would round up to 23. Always round sample size estimates up in power analyses. If we are correct that our coin lands heads 75% of the time, we need to flip at least 23 times to have an 80% chance of correctly rejecting the null hypothesis at the 0.05 significance level.
- Available tests (code chunk not evaluated):

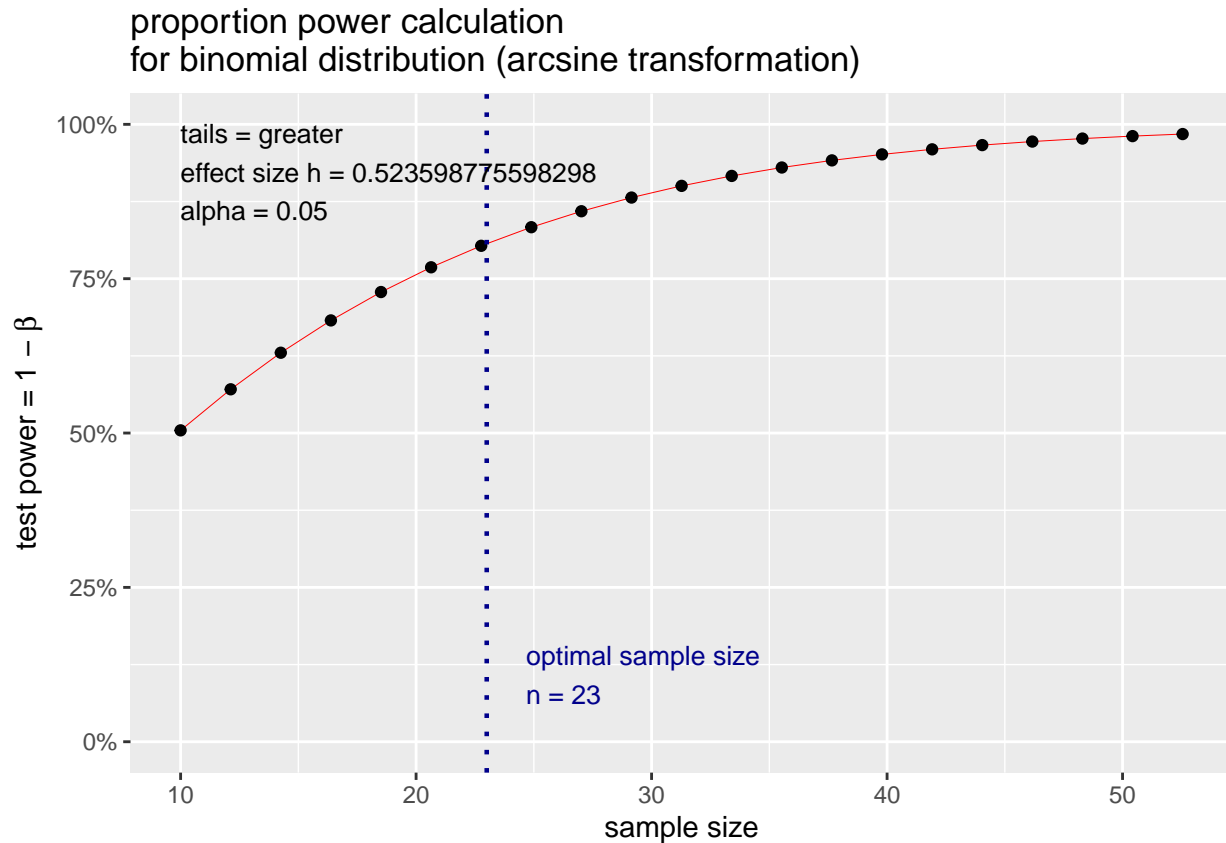
```
# two proportions (equal n)
pwr.2p.test
#two proportions(unequal n)
pwr.2p2n.test
# balanced one-way ANOVA
pwr.anova.test
# chi-square test
pwr.chisq.test
# general linear model
pwr.f2.test
# proportion (one sample)
pwr.p.test
# correlation
pwr.r.test
# t-tests (one sample, 2 sample, paired)
pwr.t.test
# t-test (two samples with unequal n)
pwr.t2n.test
```

```
pwr.p.test(h=ES.h(p1=0.75, p2=0.50),
            sig.level = 0.05,
            power = 0.80,
            alternative = "greater")
```

```
##
##      proportion power calculation for binomial distribution (arcsine transformation)
##
##              h = 0.5235988
##              n = 22.55126
##      sig.level = 0.05
##      power = 0.8
##      alternative = greater
```



```
p.out <- pwr.p.test(h = ES.h(p1 = 0.75, p2 = 0.50),
  sig.level = 0.05,
  power = 0.80,
  alternative = "greater")
plot(p.out)
```



Example: What is the power of our test if we flip the coin 40 times and lower our Type I error tolerance to 0.01? Notice we leave out the power argument, add  $n = 40$ , and change  $\text{sig.level} = 0.01$ :

The power of our test is now about 84%. If we wish to assume a “two-sided” alternative, we can simply leave it out of the function. Notice how our power estimate drops below 80% when we do this.

```
pwr.p.test(h = ES.h(p1 = 0.75, p2 = 0.50),
  sig.level = 0.01,
  n = 40,
  alternative = "greater")
```

```
##
##      proportion power calculation for binomial distribution (arcsine transformation)
##
##              h = 0.5235988
##              n = 40
##      sig.level = 0.01
##      power = 0.8377325
##      alternative = greater
pwr.p.test(h = ES.h(p1 = 0.75, p2 = 0.50),
  sig.level = 0.01,
```

```

n = 40)

##
##      proportion power calculation for binomial distribution (arcsine transformation)
##
##              h = 0.5235988
##              n = 40
##      sig.level = 0.01
##      power     = 0.7690434
##      alternative = two.sided

```

## 21 ANOVA tests (Analysis of Variance)

ANOVA is a statistical test for estimating how a quantitative dependent variable changes according to the levels of one or more categorical independent variables. ANOVA tests whether there is a difference in means of the groups at each level of the independent variable.

Example: You can check the levels of poison with the following. You should see three character values because you convert them in factor with the mutate verb.

```

PATH<- "https://raw.githubusercontent.com/guru99-edu/R-Programming/master/poisons.csv"
df <- read.csv(PATH) %>%
  select(-X) %>%
  mutate(poison = factor(poison, ordered = TRUE))
glimpse(df)

```

```

## Rows: 48
## Columns: 3
## $ time    <dbl> 0.31, 0.45, 0.46, 0.43, 0.36, 0.29, 0.40, 0.23, 0.22, 0.21, ...
## $ poison  <ord> 1, 1, 1, 1, 2, 2, 2, 2, 3, 3, 3, 3, 1, 1, 1, 1, 2, 2, 2, ...
## $ treat   <chr> "A", "A", "A", "A", "A", "A", "A", "A", "A", "A", "A", "A", ...

```

```
levels(df$poison)
```

```
## [1] "1" "2" "3"
```

- Compute the mean and the standard deviation

```

df %>%
  group_by(poison)%>%
  summarise(
    count_poison = n(),
    mean_time = mean(time, na.rm = TRUE),
    sd_time = sd(time, na.rm=TRUE)
  )

```

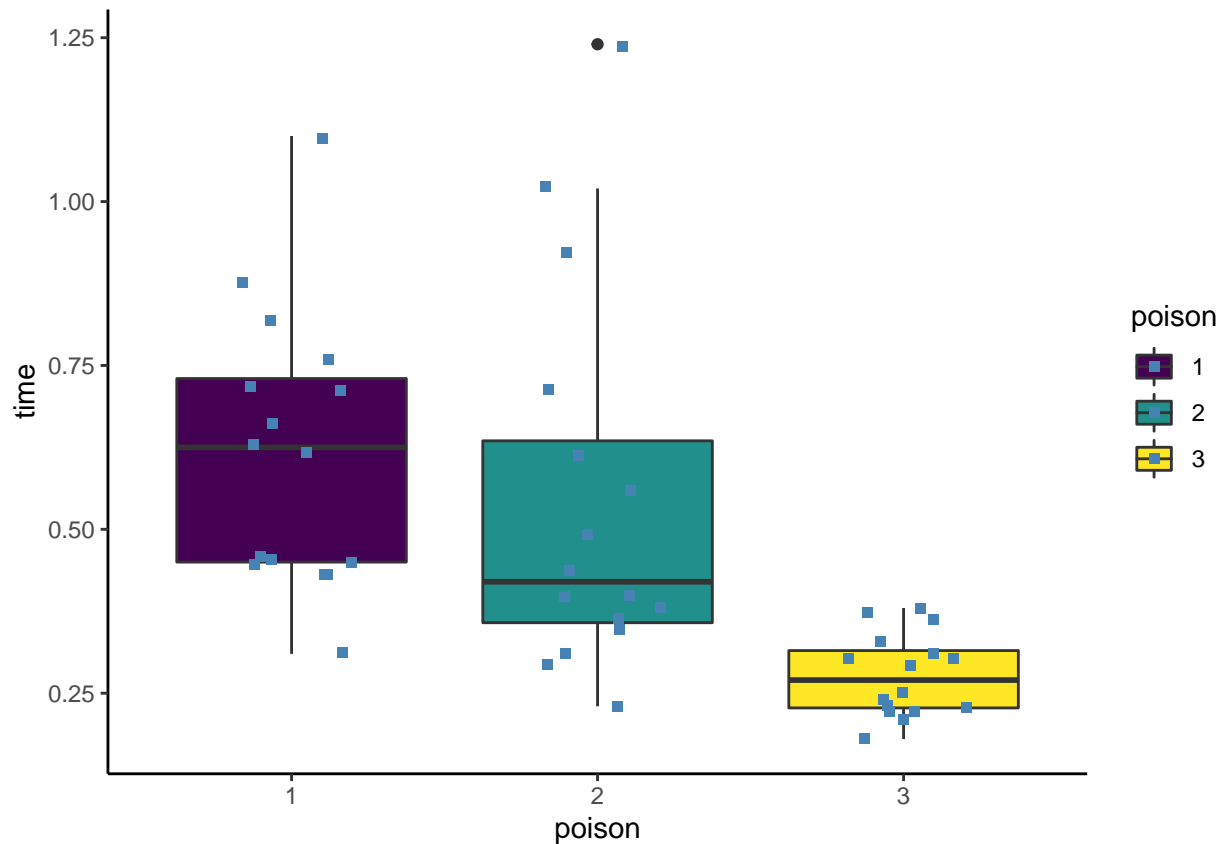
```

## # A tibble: 3 x 4
##   poison count_poison mean_time sd_time
##   <ord>         <int>     <dbl>   <dbl>
## 1 1             16      0.618   0.209
## 2 2             16      0.544   0.289
## 3 3             16      0.276   0.0623

```

- Graphically check if there is a difference between the distribution.

```
ggplot(df, aes(x = poison, y = time, fill = poison)) +
  geom_boxplot() +
  geom_jitter(shape = 15,
             color = "steelblue",
             position = position_jitter(0.21)) +
  theme_classic()
```



Run the one-way ANOVA test with the command aov.

- Syntax: `aov(formula, data)`

```
anova_one_way <- aov(time~poison, data = df)
summary(anova_one_way)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## poison      2  1.033   0.5165    11.79 7.66e-05 ***
## Residuals  45  1.972   0.0438
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

- Tukeys Test:

```
TukeyHSD(anova_one_way)
```

```
##    Tukey multiple comparisons of means
##      95% family-wise confidence level
##
## Fit: aov(formula = time ~ poison, data = df)
##
## $poison
```

```
##           diff           lwr           upr           p adj
## 2-1 -0.073125 -0.2525046  0.10625464 0.5881654
## 3-1 -0.341250 -0.5206296 -0.16187036 0.0000971
## 3-2 -0.268125 -0.4475046 -0.08874536 0.0020924
```

- ANOVA two-way test

```
anova_two_way <- aov(time~poison + treat, data = df)
summary(anova_two_way)
```

```
##           Df Sum Sq Mean Sq F value    Pr(>F)
## poison      2  1.0330   0.5165    20.64 5.7e-07 ***
## treat       3  0.9212   0.3071    12.27 6.7e-06 ***
## Residuals   42  1.0509   0.0250
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

## 22 Nonlinear Regression

- Load the data

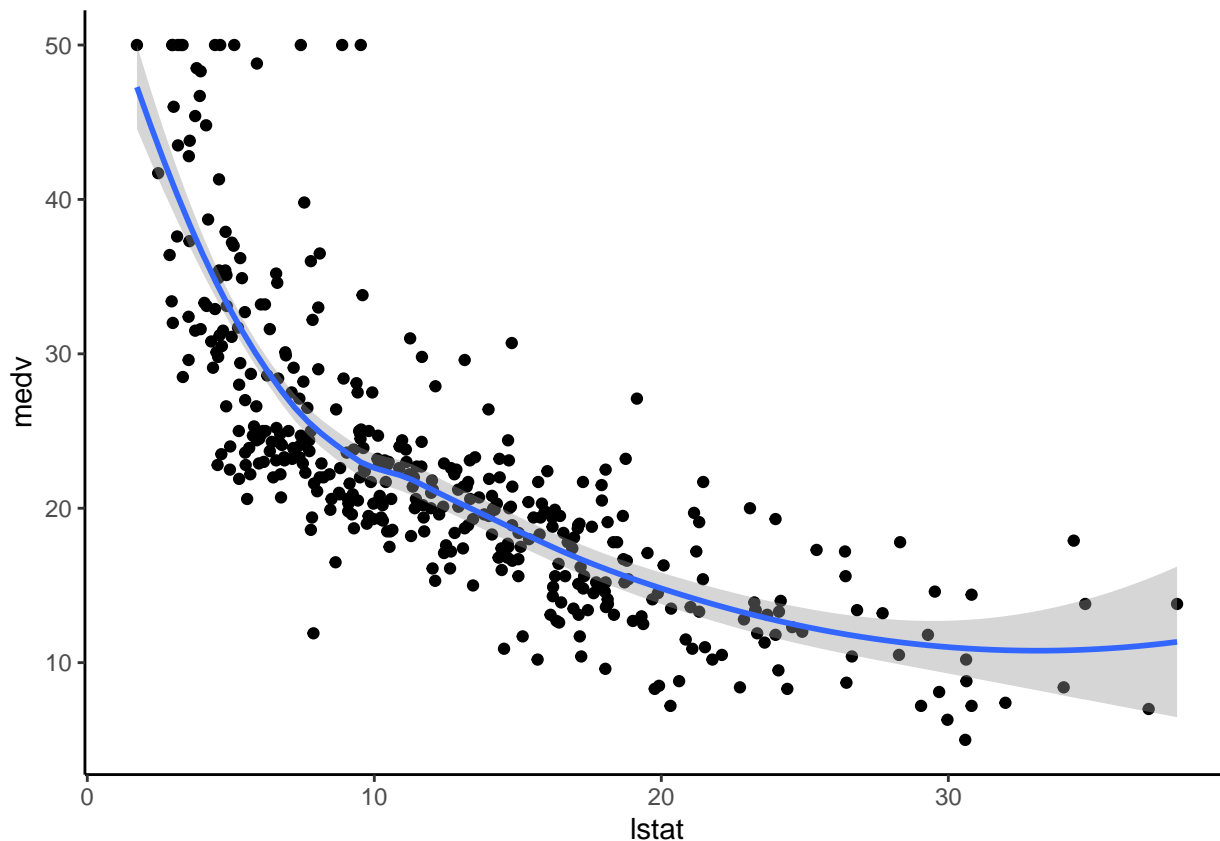
```
theme_set(theme_classic())
data("Boston", package = "MASS")
```

- Split the data into training and test set

```
set.seed(123)
training.samples <- Boston$medv %>%
  createDataPartition(p = 0.8, list = FALSE)
train.data <- Boston[training.samples, ]
test.data <- Boston[-training.samples, ]
```

- First, visualize the scatter plot of the medv vs lstat variables:

```
ggplot(train.data, aes(lstat, medv) ) +
  geom_point() +
  stat_smooth()
```



The standard linear regression model equation can be written as:

$$medv = b_0 + b_1 \times l_{stat}$$

Compute the linear regression model: \* Build the model

```
model <- lm(medv ~ lstat, data = train.data)
```

- Make predictions

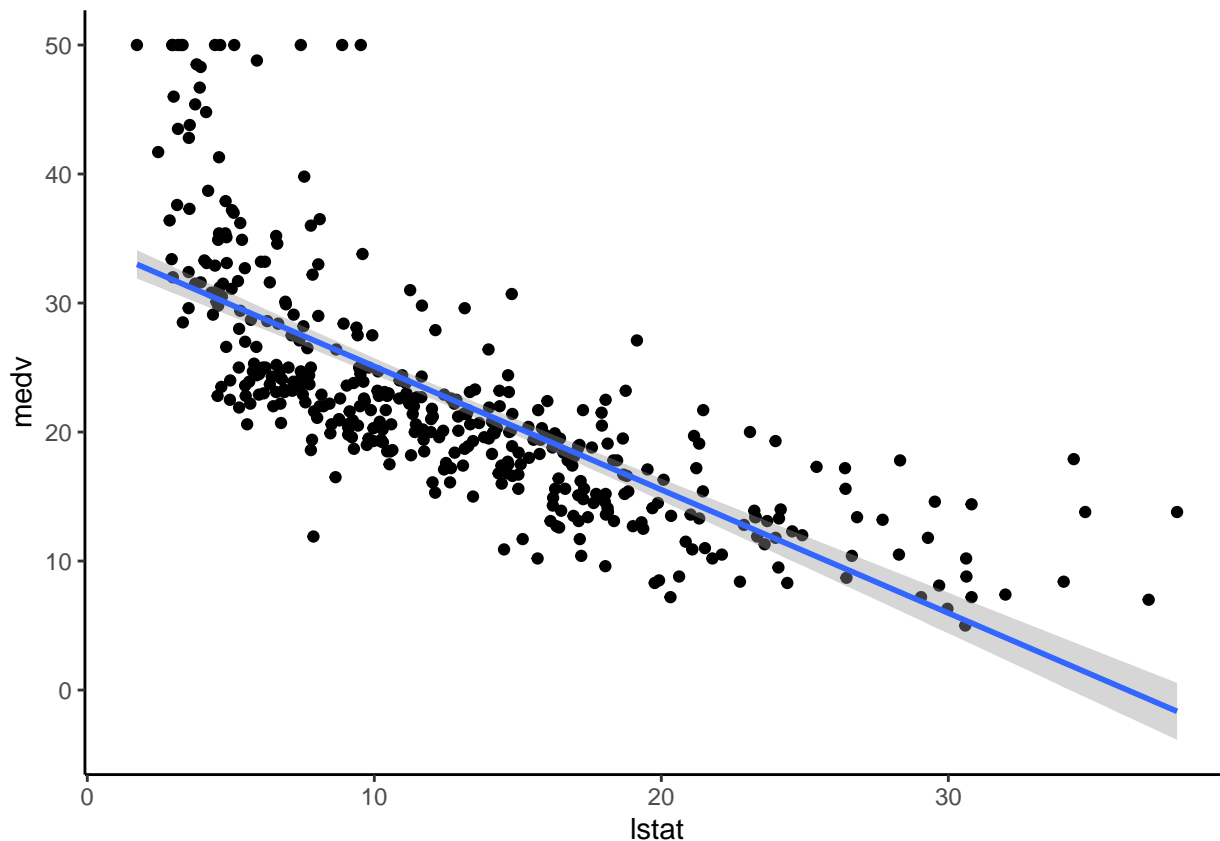
```
predictions <- model %>% predict(test.data)
```

- Model performance

```
data.frame(
  RMSE = RMSE(predictions, test.data$medv),
  R2 = R2(predictions, test.data$medv)
)
```

```
##      RMSE      R2
## 1 6.503817 0.513163
```

```
ggplot(train.data, aes(lstat, medv) ) +
  geom_point() +
  stat_smooth(method = lm, formula = y ~ x)
```



## 23 Polynomial Regression

Polynomial regression: adds polynomial or quadratic terms to the regression equations

```
lm(medv ~ lstat + I(lstat^2), data = train.data)
```

```
##
## Call:
## lm(formula = medv ~ lstat + I(lstat^2), data = train.data)
##
## Coefficients:
## (Intercept)      lstat      I(lstat^2)
##    42.5736     -2.2673      0.0412
```

```
#or
lm(medv ~ poly(lstat, 2, raw = TRUE), data = train.data)
```

```
##
## Call:
## lm(formula = medv ~ poly(lstat, 2, raw = TRUE), data = train.data)
##
## Coefficients:
##              (Intercept)  poly(lstat, 2, raw = TRUE)1
##              42.5736              -2.2673
## poly(lstat, 2, raw = TRUE)2
##              0.0412
```

```
lm(medv ~ poly(lstat, 6, raw = TRUE), data = train.data) %>%
  summary()

##
## Call:
## lm(formula = medv ~ poly(lstat, 6, raw = TRUE), data = train.data)
##
## Residuals:
```

	Min	1Q	Median	3Q	Max
##	-13.1962	-3.1527	-0.7655	2.0404	26.7661

```
##
## Coefficients:
```

		Estimate	Std. Error	t value	Pr(> t )
##	(Intercept)	7.788e+01	6.844e+00	11.379	< 2e-16 ***
##	poly(lstat, 6, raw = TRUE)1	-1.767e+01	3.569e+00	-4.952	1.08e-06 ***
##	poly(lstat, 6, raw = TRUE)2	2.417e+00	6.779e-01	3.566	0.000407 ***
##	poly(lstat, 6, raw = TRUE)3	-1.761e-01	6.105e-02	-2.885	0.004121 **
##	poly(lstat, 6, raw = TRUE)4	6.845e-03	2.799e-03	2.446	0.014883 *
##	poly(lstat, 6, raw = TRUE)5	-1.343e-04	6.290e-05	-2.136	0.033323 *
##	poly(lstat, 6, raw = TRUE)6	1.047e-06	5.481e-07	1.910	0.056910 .

```
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 5.188 on 400 degrees of freedom
## Multiple R-squared:  0.6845, Adjusted R-squared:  0.6798
## F-statistic: 144.6 on 6 and 400 DF,  p-value: < 2.2e-16
```

- Build the model

```
model <- lm(medv ~ poly(lstat, 5, raw = TRUE), data = train.data)
```

- Make predictions

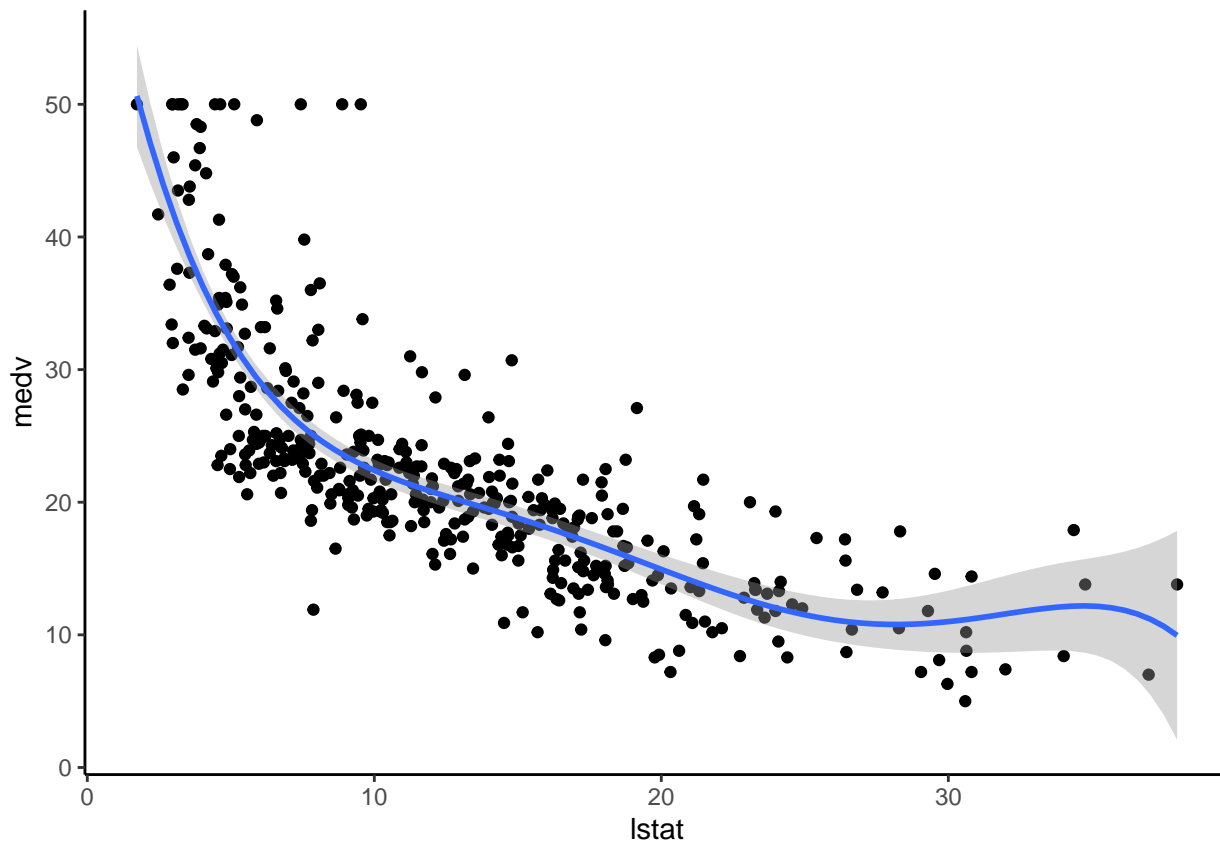
```
predictions <- model %>% predict(test.data)
```

- Model performance

```
data.frame(
  RMSE = RMSE(predictions, test.data$medv),
  R2 = R2(predictions, test.data$medv)
)
```

```
##      RMSE      R2
## 1 5.270374 0.6829474
```

```
ggplot(train.data, aes(lstat, medv)) +
  geom_point() +
  stat_smooth(method = lm, formula = y ~ poly(x, 5, raw = TRUE))
```



## 24 Log Transformation

Log transformation: when you have a non-linear relationship, you can try a logarithmic transformation of the predictor variables:

- Build the model

```
model <- lm(medv ~ log(lstat), data = train.data)
```

- Make predictions

```
predictions <- model %>% predict(test.data)
```

- Model performance

```
data.frame(
  RMSE = RMSE(predictions, test.data$medv),
  R2 = R2(predictions, test.data$medv)
)
```

```
##      RMSE      R2
## 1 5.467124 0.6570091
```

```
ggplot(train.data, aes(lstat, medv)) +
  geom_point() +
  stat_smooth(method = lm, formula = y ~ log(x))
```



