

Statistics with R

W. John Braun, UBC

COSC 405, DATA 405, and DATA 505

Winter Term 2019/20



Box plots

The `rivers` data set contains the lengths of 141 important or major North American rivers. A quick numeric summary of these data is obtained through

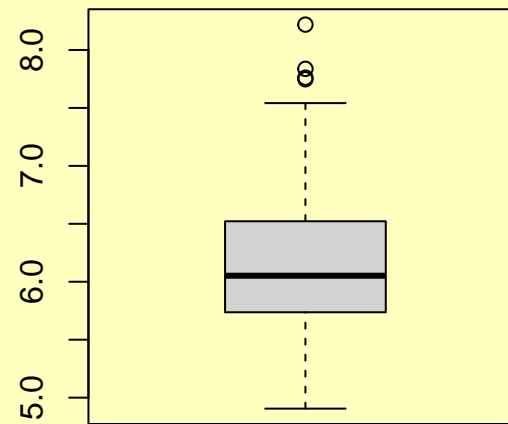
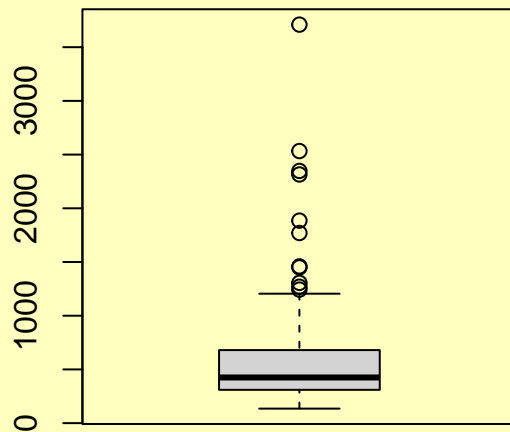
```
summary(rivers)
```

```
##      Min. 1st Qu.  Median    Mean 3rd Qu.    Max.
##      135    310    425    591    680    3710
```

Box plots

A box plot, as shown in the left panel below, can be constructed using

```
boxplot(rivers)
```



Box plots

Taking logs and then computing the box plot gives the graph in the right panel of the previous slide.

The result is much more symmetric; the histogram would be hard to distinguish from a normal distribution.

```
boxplot(log(rivers))
```

Eye Colour - Categorical Data

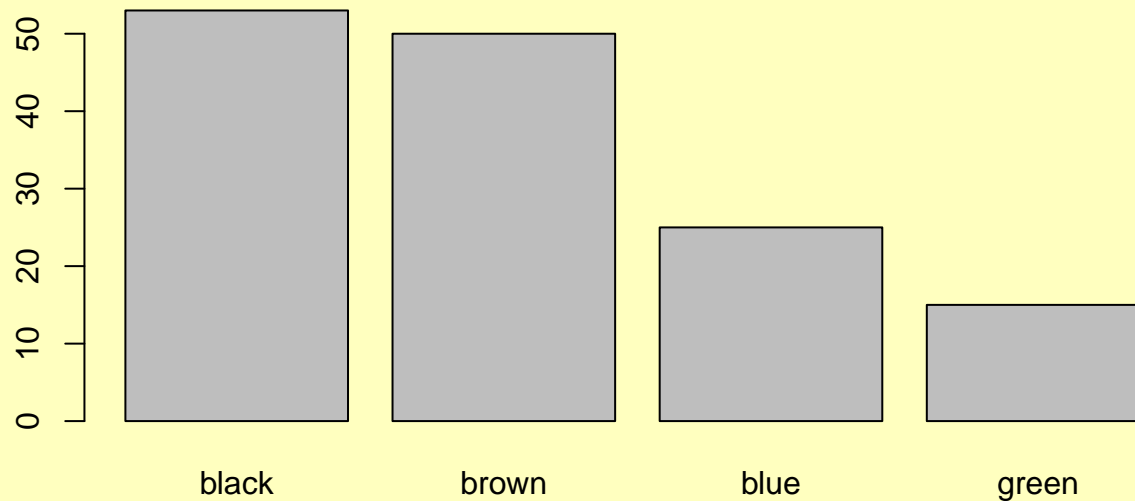
A sample of brown-haired males revealed the following eye colour counts:

black	brown	blue	green
53	50	25	15

Eye Colour - Categorical Data

A bar chart for the eye-colour counts for brown-haired males is constructed using

```
barplot(c("black" = 53, "brown" = 50, "blue" = 25, "green" = 15))
```



The goal of these tests, and the related confidence intervals, is to provide information about the mean of a single population, or about the difference in means of two population.

The critical assumption underlying the t-test is that the measurements are independent of each other.

We will use simulation to demonstrate the techniques.

One sample

We suppose that we have a random sample of measurements from a population with unknown mean μ and variance σ^2 .

Without telling you how, I will simulate such 8 such measurements, storing them in an object called `x`, and we will use a test to determine if the true mean is 0 or not:

```
## [1] 0.21 2.97 -1.37 0.90 1.36 0.86 1.44 0.46
```

We can calculate the mean and standard deviation for this sample using the `mean()` and `sd()` functions:

```
mean(x); sd(x)
```

```
## [1] 0.854
```

```
## [1] 1.23
```


One sample test

Clearly, the sample mean is not 0, but the true mean could still be 0, and this result could just be the result of random sampling error. The t-test helps us answer this question:

```
t.test(X, conf.level = .995)

##
## One Sample t-test
##
## data: X
## t = 2, df = 7, p-value = 0.09
## alternative hypothesis: true mean is not equal to 0
## 99.5 percent confidence interval:
## -0.898 2.605
## sample estimates:
## mean of x
## 0.854
```

One sample test

The small p-value indicates strong evidence against the hypothesis that the true mean is 0. In fact, this assertion is correct, since the code used to generate the random sample was as follows:

```
X <- rnorm(8, mean = 1.5) # true mean is 1.5
```

Note that we have used a 99.5% confidence interval to estimate the mean. This differs from the usual 95% that you might have been told to use.

Two samples - matched pairs

If there is a one-to-one correspondence between measurements in one of the samples with measurements in the other sample, then the appropriate way to compare the means is by taking the differences, and running a one-sample test on the differences.

This can be done with the `paired` option in the `t.test()` function.

Let's suppose L is a set of left foot lengths (in cm) for a sample of 15 adult males, and R contains the corresponding right foot lengths.

We would be interested in any systematic difference in the lengths of the feet.

Two samples - matched pairs

A simulation model for the case where there is no difference could be the following:

```
L <- rnorm(15, mean = 28, sd = 1)
R <- L + rnorm(15, mean = 0, sd = .03)
```

Here we have assumed that the left feet are normally distributed with a mean of 28 cm and a standard deviation of 1 cm.

The right feet lengths are not exactly equal to the left feet lengths, but on average there is no difference.

The standard deviation of the difference is small.

Two samples - matched pairs

Let's see what the t-test says:

```
t.test(L, R, paired=TRUE)

##
## Paired t-test
##
## data: L and R
## t = 0.5, df = 14, p-value = 0.6
## alternative hypothesis: true mean difference is not equal to 0
## 95 percent confidence interval:
## -0.0121 0.0191
## sample estimates:
## mean difference
## 0.00345
```

The p-value is large indicating that there is no evidence of a difference, in line with the truth.

Simple Regression

The yield (y , in kg/plot) was measured for various salinity concentrations (x , measured in units of electrical conductivity).

18 measurements were recorded in a file called *tomato.txt* whose contents summarized below:

```
##      salinity      yield
##  Min.      : 1.60    Min.      :41.0
##  1st Qu.: 2.15    1st Qu.:48.6
##  Median : 4.90    Median :53.0
##  Mean    : 5.46    Mean    :52.4
##  3rd Qu.: 9.15    3rd Qu.:56.4
##  Max.    :10.20    Max.    :63.1
```

The first column contains the salinity concentration levels, and the second column contains the yield measurements.

Simple Regression

We read these data into R using the `read.table()` function (or using a menu option in RStudio):

```
tom <- read.table("tomato.txt", header=FALSE)
```

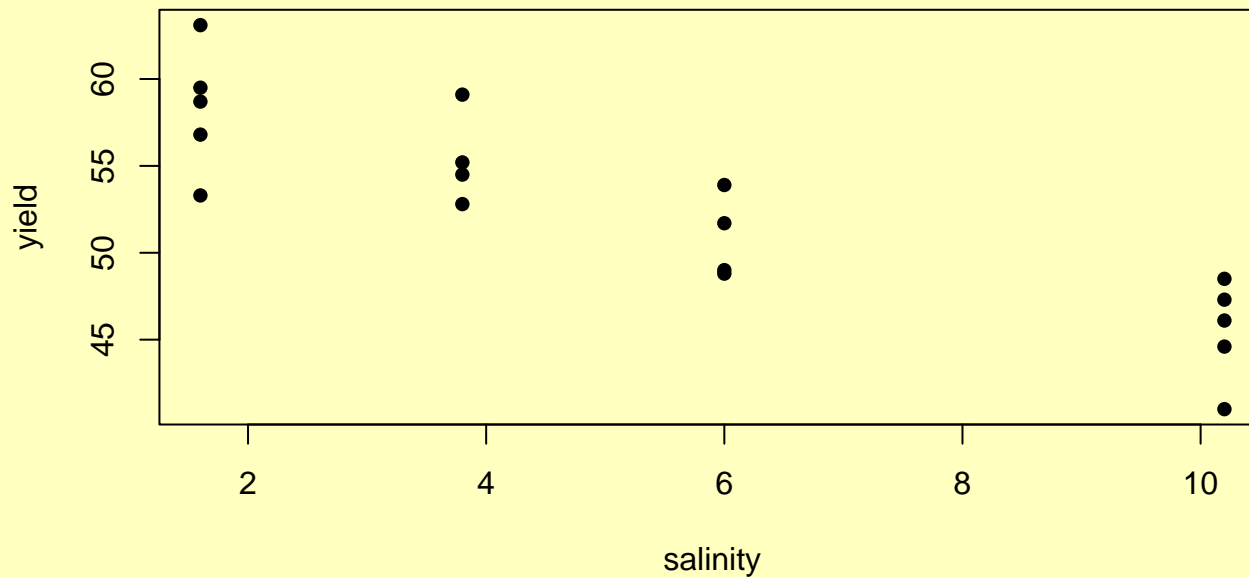
Since there is no header, we should apply some sensible names to the data frame:

```
names(tom) <- c("salinity", "yield")
```

Simple Regression

We next construct a *scatterplot* of the data to look for patterns and outliers.

```
plot(yield ~ salinity, data = tom, pch=16)
```



Simple Regression

The scatterplot gives an indication of a clear downward trend as salinity increases.

The trend is also vaguely linear.

We can investigate this with the `lm()` function:

```
tom.lm <- lm(yield ~ salinity, data = tom)
```

Simple Regression

We can explore the output from the fitted model using the `summary()` function:

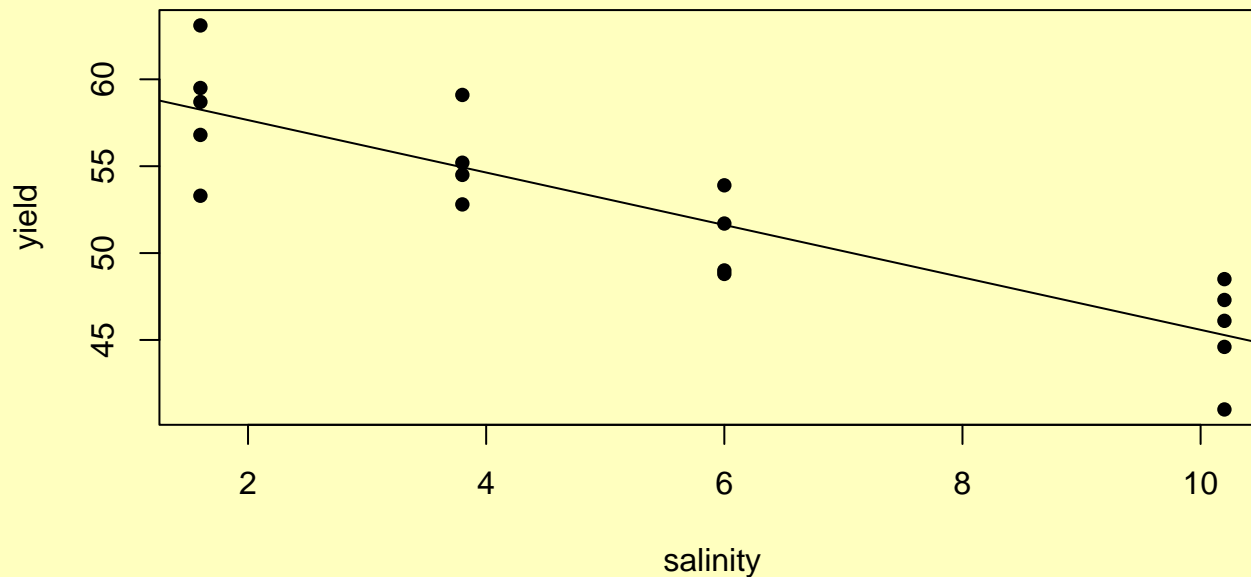
```
summary(tom.lm)

##
## Call:
## lm(formula = yield ~ salinity, data = tom)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -4.956  -1.967   0.173   1.825   4.844
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)    60.67      1.28    47.37 < 2e-16
## salinity       -1.51      0.20   -7.53 1.2e-06
##
## Residual standard error: 2.83 on 16 degrees of freedom
## Multiple R-squared:  0.78, Adjusted R-squared:  0.766
## F-statistic: 56.7 on 1 and 16 DF,  p-value: 1.21e-06
```

Simple Regression

We can overlay the scatterplot of the data with the fitted line using the `abline()` function and the output from the `lm()` function:

```
abline(tom.lm)
```



One factor ANOVA

The `chickwts` data frame contains measurements of the weights of chicks who have been randomly assigned to groups, each of which has been given a different type of feed.

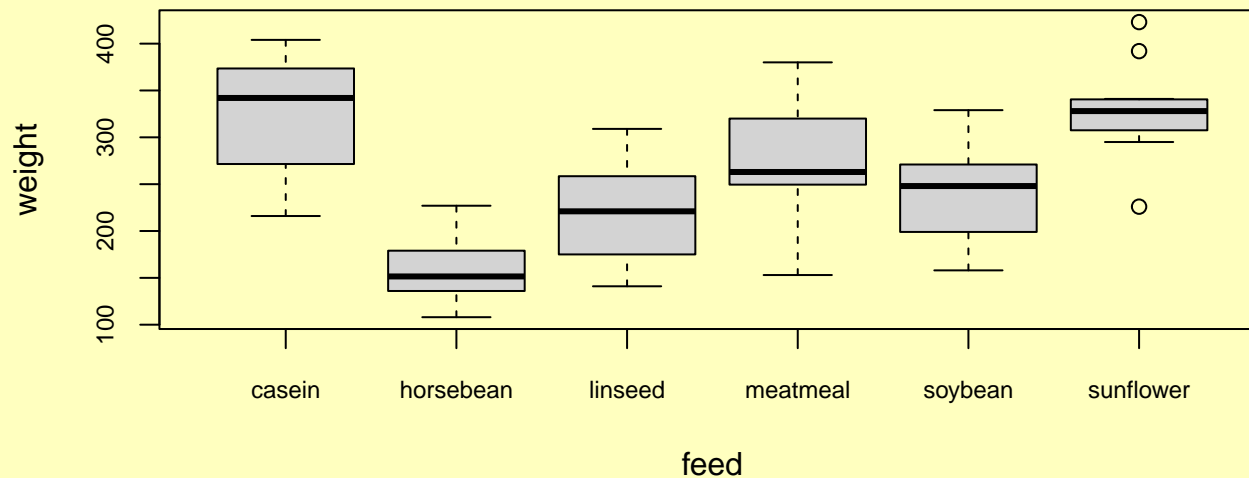
It is of interest to know whether the different feed types lead to systematic differences in weight.

We refer to feed type as a *factor* having different *levels* representing the particular kinds of feed, e.g. linseed, horsebean, and so on.

One factor ANOVA

Side-by-side *boxplots*, as displayed below, are a useful way to visualize these data.

```
plot(weight ~ feed, data = chickwts, cex.axis=.75)
```



Because `feed` is a factor, `plot()` yields box plots.

One factor ANOVA

From the graph, it seems that horsebean leads to lower weights than some of the other feed types.

It is hard to tell for sure if there is variability between treatments because of the variability within treatments, that is noise due to unmeasured factors.

We can test whether there is a difference in the mean weights statistically with the analysis of variance (ANOVA).

One factor ANOVA

A general purpose procedure is as follows:

```
chick.lm <- lm(weight ~ feed, data = chickwts)
anova(chick.lm)

## Analysis of Variance Table
##
## Response: weight
##           Df Sum Sq Mean Sq F value Pr(>F)
## feed       5 231129   46226   15.4 5.9e-10
## Residuals 65 195556    3009
```

One factor ANOVA

The test statistic compares the variability in the averages with the variability in the noise through an F -statistic.

A p -value is computed which gives the strength of evidence against the *null hypothesis*, that is the hypothesis that there is no difference in the means.

A small p -value – and in this case, it is very small – indicates strong evidence against the null hypothesis, in favour of the alternative that there is a difference.

Multiple Regression

The data frame `table.b4` in the *MPV* library contains the following columns:

```
y sale price of the house (in thousands of dollars)
x1 taxes (in thousands of dollars)
x2 number of baths
x3 lot size (in thousands of square feet)
x4 living space (in thousands of square feet)
x5 number of garage stalls
x6 number of rooms
x7 number of bedrooms
x8 age of the home (in years)
x9 number of fireplaces
```

Multiple Regression

There are 24 observations on these *variables* in the data frame.

A natural question to ask is whether any or all of the given variables or *covariates* could be used to predict the sale price of a house.

We consider a linear model of the form

$$y = \beta_0 + \sum_{j=1}^9 \beta_j x_j + \varepsilon.$$

The elements of ε are assumed to be uncorrelated random variables with mean 0 and common variance σ^2 .

Fitting the model

The `lm()` function will take care of the coefficient estimation, variance estimation, t and F and p -value calculations in one function call.

For example, if we want to relate house price, y to x_1 and x_4 ,

```
house.lm <- lm(y ~ x1 + x4, data=table.b4)
```

We can view the output from this, using the `summary` function as in

```
summary(house.lm)$coefficients
```

	Estimate	Std. Error	t value	Pr(> t)
## (Intercept)	11.54	3.175	3.636	0.001544
## x1	2.92	0.576	5.069	0.000051
## x4	3.16	3.298	0.957	0.349315

Estimating and predicting

The model can now be used to estimate the expected house price for houses with x_1 taxes and x_4 amount of living space using the formula

$$\hat{y} = 11.5 + 2.92x_1 + 3.15x_4.$$

This can be accomplished in R using the `predict()` function.

For instance, suppose we want to estimate the mean sale price for homes with \$2000 taxes and 3000 square feet of living area. Use

```
predict(house.lm, newdata = data.frame(x1 = 2, x4 = 3))

##      1
## 26.9
```

Estimating and predicting

For interval estimation, use `interval = "confidence"` **for estimation and** `interval = "predict"` **for prediction. For interval estimation of a prediction:**

```
predict(house.lm, newdata = data.frame(x1 = 2, x4 = 3),
        interval = "predict")
```

```
##      fit  lwr  upr
## 1 26.9 10.2 43.5
```

Interval estimation for the mean:

```
predict(house.lm, newdata = data.frame(x1 = 2, x4 = 3),
        interval = "confidence")
```

```
##      fit  lwr  upr
## 1 26.9 11.4 42.3
```

The *analysis of covariance* (ANCOVA) allows us to model continuous responses as linear functions of continuous and categorical covariates.

It can be viewed as a relatively straightforward extension of multiple regression.

It can also be viewed as an extension of ANOVA whereby there is a blocking factor which is continuously measured.

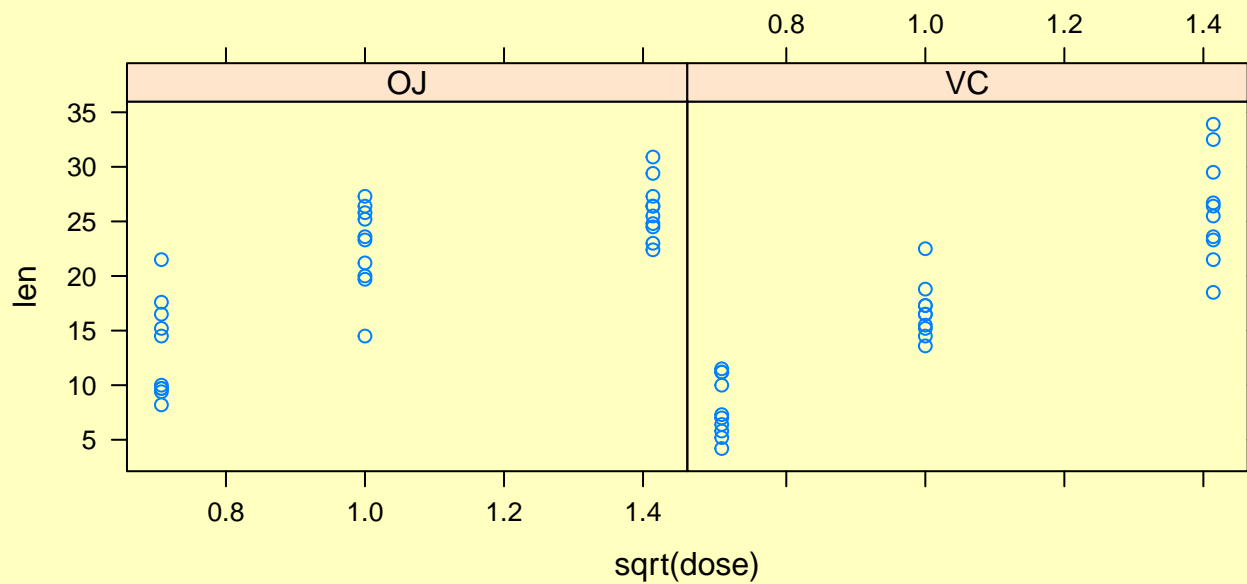
For a categorical covariate with two levels, there would be two lines in the regression: parallel if there is no interaction effect; two different slopes if there is an interaction effect.

The `ToothGrowth` data frame in R concerns the length of odontoblasts, cells connected with the growth of teeth, in a sample of 60 guinea pigs. One of three dose levels of vitamin C were supplied to the guinea pigs in one of two forms: `supp = VC` refers to ascorbic acid and `supp = OJ` refers to orange juice.

The data are displayed on the next slide, using the `xypplot()` function from the `lattice` package (Sarkar, 2008).

ANCOVA

```
library(lattice)
xyplot(len ~ sqrt(dose) | supp, data = ToothGrowth)
```



The figure shows that there are possibly two different lines relating length to vitamin C dose; it is possible that there is a treatment effect.

We use `lm()` to check this, first by allowing for two intercepts and two slopes:

```
TG.lm <- lm(len ~ sqrt(dose)*supp, data = ToothGrowth)
summary(TG.lm)$coefficients
```

##	Estimate	Std. Error	t value	Pr(> t)
## (Intercept)	2.48	2.63	0.943	3.50e-01
## sqrt(dose)	17.48	2.43	7.185	1.71e-09
## suppVC	-12.02	3.72	-3.236	2.04e-03
## sqrt(dose):suppVC	8.00	3.44	2.325	2.37e-02

Looking at the interaction between the square root of dose and supp , we see a fairly small p -value which is highly suggestive of different slopes.

There is no reason to consider the model without different slopes (which would have been obtained by replacing the $*$ with $+$). The value 8 indicates that for VC , the slope is 8 units higher than for OJ :

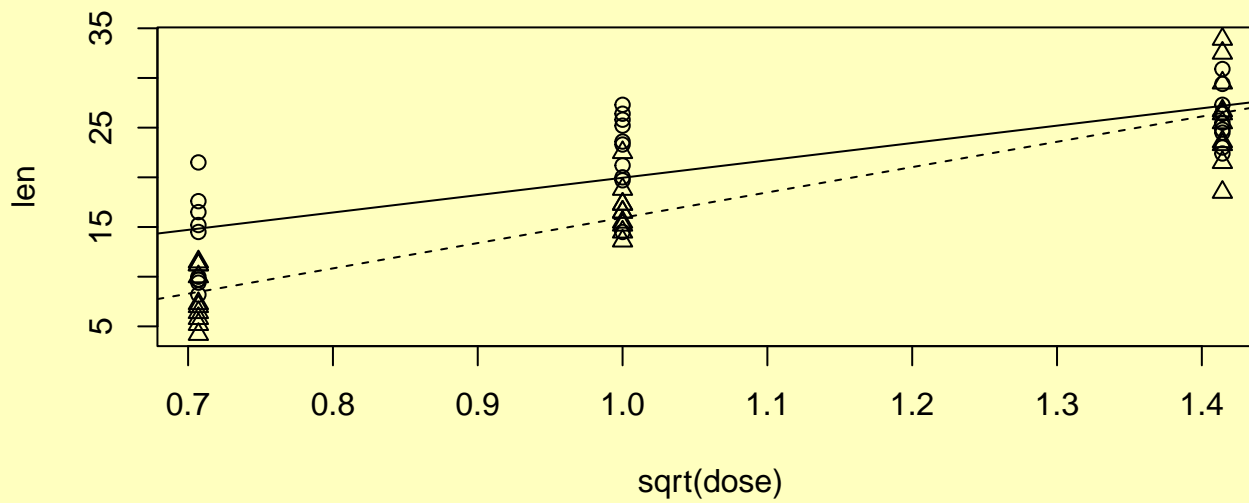
$$17.48 + 8 = 25.48.$$

The intercept for VC is 12.02 units lower: $2.478 - 12.024 = -9.546$.

ANCOVA

We can graphically summarize the data with a simple scatterplot with the lines overlaid:

```
plot(len ~ sqrt(dose), pch = as.numeric(supp),
     data = ToothGrowth)
abline(2.478, 17.479) # OJ line
abline(2.478-12.024, 17.479 + 8, lty=2) # VC line, dashed
```



Modelling binary responses with logistic regression

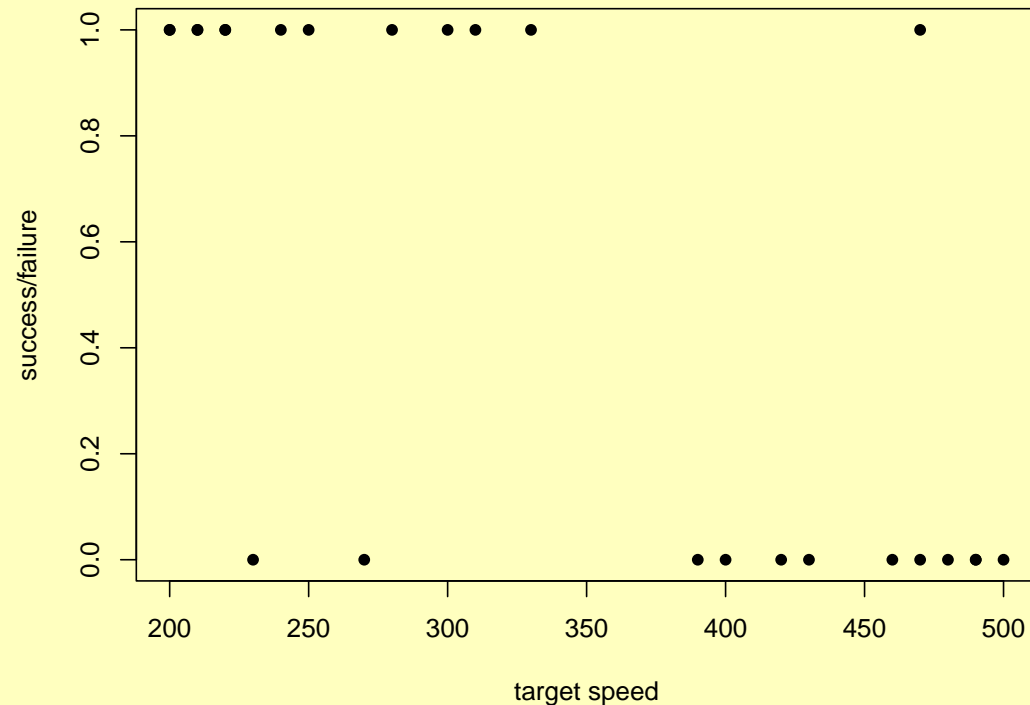
The data in `p13.1` in the *MPV* package describes successes and failures of surface-to-air missiles as they relate to target speed.

Such *binary data* are not nearly normally distributed, so the efficacy of least-squares becomes very questionable here.

We now suggest some of the R functions that can be used to begin to analyze such data.

Modelling binary responses with logistic regression

The data are plotted to the right, with successes on the vertical axis being represented by a '1' and failures being represented by a '0'.



```
library(MPV) # contains the p13.1 data frame
plot(p13.1, xlab = "target speed", ylab = "success/failure",
      pch=16) # produces a different plotting character
```

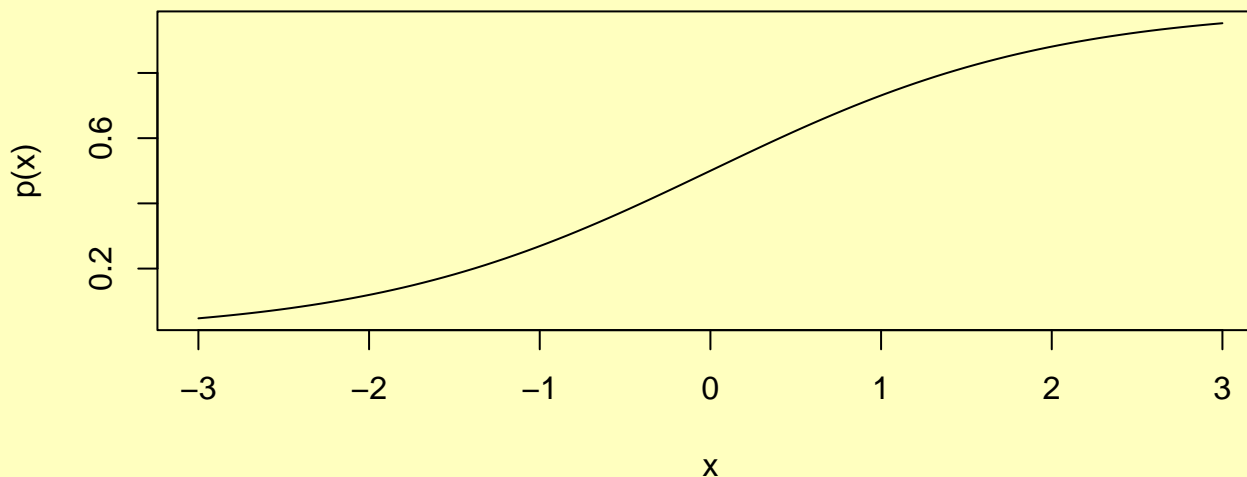
Modelling binary responses with logistic regression

Perhaps the most popular function for this purpose is the *logistic* function

$$p(x) = \frac{e^x}{e^x + 1}.$$

The function is sketched below.

```
curve(exp(x) / (1 + exp(x)), from = -3, to = 3, ylab="p(x)")
```



Modelling binary responses with logistic regression

The inverse of the logistic function is logit function:

$$\ell(p) = \log \left(\frac{p}{1-p} \right).$$

While p is restricted to take values between 0 and 1, the logit function can take any possible value, so relating the logit function to a straight line or other linear combination is a possibility. That is,

$$\ell(p(x)) = \beta_0 + \beta_1 x$$

which means that we can express the probability of an event in terms of a covariate x .

$\ell()$ is an example of a *link* function.

Modelling binary responses with logistic regression

To fit the logistic regression model to the missile success data, try

```
p13.glm <- glm(y ~ x, data = p13.1, family = binomial)
```

To see the results, try

```
summary(p13.glm)
```

Note that we did not specify the link function; the default choice with the binomial family is the logit.

Modelling binary responses with logistic regression

If we only want to see the coefficients and their standard errors, try

```
summary(p13.glm)$coefficients
```

```
##              Estimate Std. Error z value Pr(>|z|)
## (Intercept)    6.0709    2.10900    2.88  0.00399
## x              -0.0177    0.00608   -2.91  0.00357
```

This output tells us that the logit of the probability of success as a linear function of target speed has intercept 6.07 and slope -.0177.

Standard error estimates for these parameter estimates are supplied and indicate, in particular, that the slope is clearly negative.

Modelling binary responses with logistic regression

The logistic curve can now be used to calculate probabilities of success at the various speeds.

We will use the `predict()` function with `type = "response"`.

If `type` is not specified, the default is to use the predictions on the linear scale.

Modelling binary responses with logistic regression

```
plot(p13.1, pch=16, xlab = "target speed",
     ylab = "success/failure")
newspeeds <- 200:500 # speeds at which we can predict
lines(newspeeds, predict(p13.glm, newdata=
  data.frame(x = newspeeds), type = "response"))
```

