

R-Adelaide 2016

Day 3:
Linear and generalised linear models in R

Dr Steven Delean - Benham Rm 103
steven.delean@adelaide.edu.au

School of Biological Sciences - The University of Adelaide



COMMONWEALTH OF AUSTRALIA
Copyright Regulations 1969

WARNING

This material has been reproduced and communicated to you by or on behalf of The University of Adelaide under Part VII of the Copyright Act 1968 (the Act).

The material in this communication may be subject to copyright under the Act. Any further reproduction or communication of this material by you may be the subject of copyright protection under the Act.

Do not remove this notice.

Linear model design matrix

- Two R functions, `formula` and `model.matrix` are required to produce *design matrices* for a variety of models

$$Y_i = \beta_0 + \beta_1 x_i + \varepsilon_i, \quad i = 1, \dots, N$$

- Focus on the *observational unit* or *experimental unit* to N different entities from which we obtain a measurement
- For categorical variables often called *dummy* or *indicator variables*
 - Simply indicate if experimental unit had a certain characteristic or not

Linear model design matrix

Using linear model algebra to represent model for design matrix \mathbf{X} :

$$\mathbf{Y} = \begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_N \end{pmatrix}, \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix}, \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \text{ and } \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{pmatrix}$$

as

$$\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_N \end{pmatrix} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_N \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_N \end{pmatrix}$$

or simply

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

- Define a design matrix \rightarrow find the least squares estimates by *fitting the model*
- In R, we directly provide a *formula* to the model fitting function (e.g., `lm`)
- The formula statement is central to model fitting in R generally
- Generates a `model.matrix`, which is used internally by `lm` function

- Design matrix is critical step in linear modeling
 - Encodes coefficients to be fit in model, and
 - Inter-relationship between samples
- Encodes assumptions about the relationships between the response and explanatory variables
- First component usually an *intercept* column (column of 1's)
- To compare differences between groups, second column specifies which samples are in second group

The basic formula - comparing groups

- Encode experimental design with:
 - Formula with the tilde symbol '~'
 - So model observations using variables to right of tilde
 - Then specify name of a variable that classifies samples into groups
- Two groups, control and high fat diet, with two samples each
- Tell R that these values should not be interpreted numerically, but as different levels of a **factor**

```
1 # Group is a factor
2 group <- factor( c(1,1,2,2) )
3 model.matrix(~ group)
4 # Do not need to use formula function as '~' implies this
5 model.matrix(formula(~ group))
```

```
1 (Intercept) group2
2 1          1      0
3 2          1      0
4 3          1      1
5 4          1      1
6 attr(,"assign")
7 [1] 0 1
8 attr(,"contrasts")
9 attr(,"contrasts")$group
10 [1] "contr.treatment"
```

The basic formula - comparing groups

- What happens if we don't tell R that 'group' should be interpreted as a factor?
 - Provided numeric variable so Formula with tilde symbol '~'
 - Want second column to have only 0 and 1, indicating group membership
- Note: names of factor levels are irrelevant to `model.matrix` and `lm`
 - All that matters is order
- Tell R that these values should not be interpreted numerically, but as different levels of a **factor**

```
1 # Group is not a factor
2 group <- c(1,1,2,2)
3 model.matrix(~ group)
4 # Using different names produces same model.matrix
5 group <- factor(c("control","control","highfat","highfat"))
6 model.matrix(~ group)
```

```
1 (Intercept) group
2 1 1 1
3 2 1 1
4 3 1 2
5 4 1 2
6 attr(,"assign")
7 [1] 0 1
8
9 (Intercept) grouphighfat
10 1 1 0
11 2 1 0
12 3 1 1
13 4 1 1
14 attr(,"assign")
15 [1] 0 1
16 attr(,"contrasts")
17 attr(,"contrasts")$group
18 [1] "contr.treatment"
```

Formula extensions

```
1 # More groups
2 group <- factor(c(1,1,2,2,3,3))
3 model.matrix(~ group)
```

```
1 # Alternate formulation
2 group <- factor(c(1,1,2,2,3,3))
3 model.matrix(~ group + 0)
```

```
1      (Intercept) group2 group3
2      1          1      0      0
3      2          1      0      0
4      3          1      1      0
5      4          1      1      0
6      5          1      0      1
7      6          1      0      1
8 attr(,"assign")
9 [1] 0 1 1
10 attr(,"contrasts")
11 attr(,"contrasts")$group
12 [1] "contr.treatment"
```

```
1      group1 group2 group3
2      1      1      0      0
3      2      1      0      0
4      3      0      1      0
5      4      0      1      0
6      5      0      0      1
7      6      0      0      1
8 attr(,"assign")
9 [1] 1 1 1
10 attr(,"contrasts")
11 attr(,"contrasts")$group
12 [1] "contr.treatment"
```


More variables

Assume diet effect is same across sexes

$$Y_i = \beta_0 + \beta_1 x_{i,1} + \beta_2 x_{i,2} + \varepsilon_i$$

```
1 # Effect of diet and difference in sexes
2 diet <- factor(c(1,1,1,1,2,2,2,2))
3 sex <- factor(c("f","f","m","m","f","f","m","m"))
4 table(diet,sex)
5 # To fit the additive model
6 diet <- factor(c(1,1,1,1,2,2,2,2))
7 sex <- factor(c("f","f","m","m","f","f","m","m"))
8 model.matrix(~ diet + sex)
```

```
1      (Intercept) diet2 sexm
2 1             1      0      0
3 2             1      0      0
4 3             1      0      1
5 4             1      0      1
6 5             1      1      0
7 6             1      1      0
8 7             1      1      1
9 8             1      1      1
10 attr(,"assign")
11 [1] 0 1 2
12 attr(,"contrasts")
13 attr(,"contrasts")$diet
14 [1] "contr.treatment"
15
16 attr(,"contrasts")$sex
17 [1] "contr.treatment"
```

```
1 model.matrix(~ diet + sex + diet:sex)
```

```
1 model.matrix(~ diet*sex)
```

```
1      (Intercept) diet2 sexm diet2:sexm
2 1             1      0      0          0
3 2             1      0      0          0
4 3             1      0      1          0
5 4             1      0      1          0
6 5             1      1      0          0
7 6             1      1      0          0
8 7             1      1      1          1
9 8             1      1      1          1
10 attr(,"assign")
11 [1] 0 1 2 3
12 attr(,"contrasts")
13 attr(,"contrasts")$diet
14 [1] "contr.treatment"
15
16 attr(,"contrasts")$sex
17 [1] "contr.treatment"
```

Reference level, data and continuous predictors

```
1 # Releveling
2 group <- factor(c(1,1,2,2))
3 group <- relevel(group, "2")
4 model.matrix(~ group)
5 # Providing factor levels
6 group <- factor(group, levels=c("1", "2"))
7 model.matrix(~ group)
```

```
1      (Intercept) group1
2      1          1      1
3      2          1      1
4      3          1      0
5      4          1      0
6 attr(,"assign")
7 [1] 0 1
8 attr(,"contrasts")
9 attr(,"contrasts")$group
10 [1] "contr.treatment"
```

```
13      (Intercept) group2
14      1          1      0
15      2          1      0
16      3          1      1
17      4          1      1
18 attr(,"assign")
19 [1] 0 1
20 attr(,"contrasts")
21 attr(,"contrasts")$group
22 [1] "contr.treatment"
```

```
1 # Finding data
2 group <- 1:4
3 model.matrix(~ group, data=data.frame(group=5:8))
4 # Continuous variables
5 tt <- seq(0,3.4,len=4)
6 model.matrix(~ tt + I(tt^2))
7 # Different contrasts
8 d <- data.frame(time=factor(1:4))
9 model.matrix(~time,data=d, contrasts.arg= list(time="contr.sum"))
```

```
1      (Intercept) group
2      1          1      5
3      2          1      6
4      3          1      7
5      4          1      8
6 attr(,"assign")
7 [1] 0 1
8
9      (Intercept)      tt      I(tt^2)
10      1          1 0.000000000 0.000000000
11      2          1 1.133333333 1.284444444
12      3          1 2.266666667 5.137777778
13      4          1 3.400000000 11.560000000
14 attr(,"assign")
15 [1] 0 1 2
```

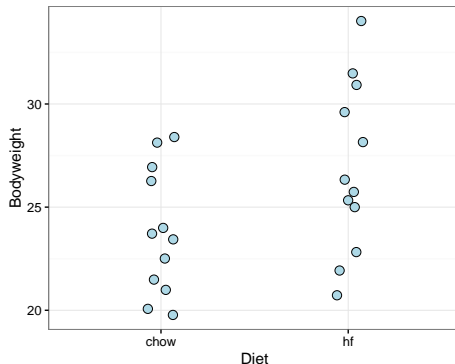
```
16
17      (Intercept) time1 time2 time3
18      1          1      1      0      0
19      2          1      0      1      0
20      3          1      0      0      1
21      4          1     -1     -1     -1
22 attr(,"assign")
23 [1] 0 1 1 1
24 attr(,"contrasts")
25 attr(,"contrasts")$group
```

Linear models - a first taste

- Does high fat diet lead to higher average body weight in mice?
 - Simple t -test problem but will demonstrate equivalence with `lm`
- High fat diet group has higher weights, but with overlap
- Build design matrix **X** using formula `~Diet`
 - Group with the 1's in second column is determined by second level of 'Diet'

```
1 dat <- read.csv("./data-raw/femaleMiceWeights.csv")
2 set.seed(1) #same jitter in stripchart
3 # Plot
4 ggplot(dat, aes(x=Diet, y=Bodyweight)) + geom_jitter(position=position_
  ↳ jitter(0.2), size=3, shape=21, bg="lightblue") + theme_bw()
5 # Design
6 levels(dat$Diet)
7 X <- model.matrix(~ Diet, data=dat)
8 head(X)
```

```
[1] "chow" "hf"
      (Intercept) Diethf
1             1      0
2             1      0
3             1      0
4             1      0
5             1      0
6             1      0
```



The Mathematics Behind `lm()`

- Quick look behind the scenes of function `lm`
- Form the design matrix \mathbf{X} and calculate β
 - Minimizes sum of squares
- Formula for solution is:

$$\hat{\beta} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{Y}$$

- Calculate in R using:
 - Matrix multiplication operator `%*%`
 - Inverse function `solve`, and
 - Transpose function `t`

```
1 Y <- dat$Bodyweight
2 X <- model.matrix(~ Diet, data=dat)
3 solve(t(X) %*% X) %*% t(X) %*% Y
```

```
1           [,1]
2 (Intercept) 23.81333333
3 Diethf      3.02083333
```

Coefficients are average of control group and difference of averages

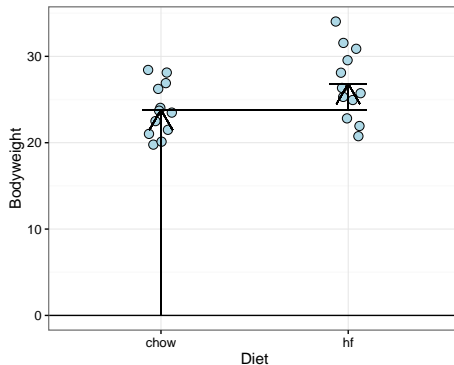
```
1 s <- with(dat, split(Bodyweight, Diet))
2 mean(s[["chow"]])
3 mean(s[["hf"]]) - mean(s[["chow"]])
```

```
1 [1] 23.81333333
2 [1] 3.02083333
```

Linear model using lm()

```
1 fit <- lm(Bodyweight ~ Diet, data=dat)
2 summary(fit)
3 (coefs <- coef(fit))
```

```
1
2 Call:
3 lm(formula = Bodyweight ~ Diet, data = dat)
4
5 Residuals:
6      Min       1Q   Median       3Q      Max
7  -6.104  -2.436  -0.414   2.834   7.186
8
9 Coefficients:
10      Estimate Std. Error t value Pr(>|t|)
11 (Intercept)  23.81      1.04   22.91  <2e-16 ***
12 Diethf       3.02      1.47    2.06   0.052 .
13 ---
14 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15
16 Residual standard error: 3.6 on 22 degrees of freedom
17 Multiple R-squared:  0.161,    Adjusted R-squared:  0.123
18 F-statistic: 4.22 on 1 and 22 DF, p-value: 0.0519
19
20 (Intercept)      Diethf
21 23.813333333  3.020833333
```

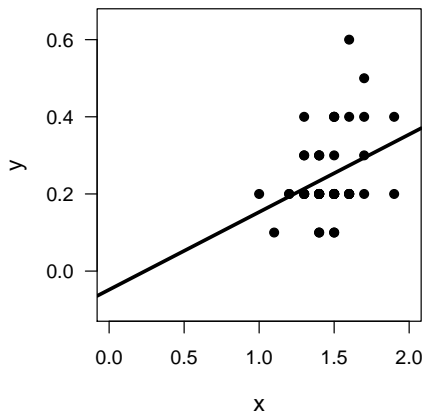


```
1 mytt <- t.test(s[["hf"]], s[["chow"]], var.equal=TRUE)
2 summary(fit)$coefficients[2, 3]
3 mytt$statistic
```

```
1 [1] 2.055173992
2      t
3 2.055173992
```

Overarching regression framework

$$y_i = \beta_0 + \beta_1 x_i + \varepsilon_i$$
$$\varepsilon_i \sim N(0, \sigma^2)$$
(1)



Data

y = response variable

x = predictor

Parameters

β_0 = intercept

β_1 = slope

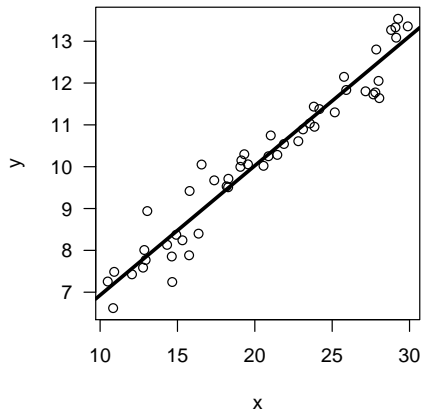
σ = residual variation

ε = residuals

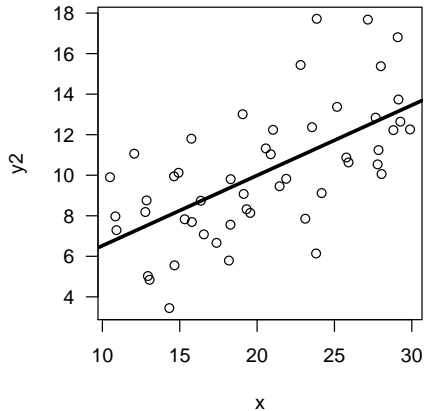
Residual variation (error)

- Our model doesn't perfectly predict Y ; residuals are used to estimate error of model

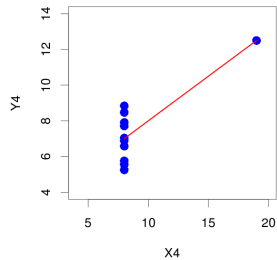
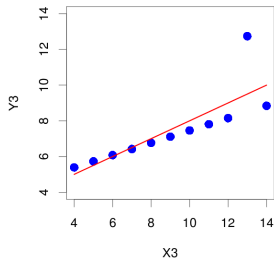
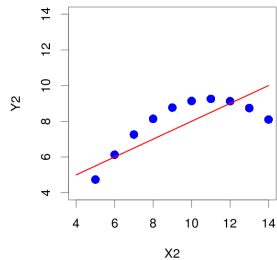
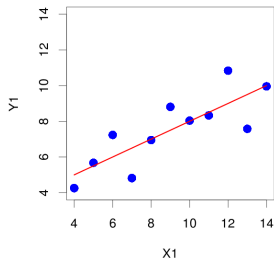
small



large



Always plot your data first!



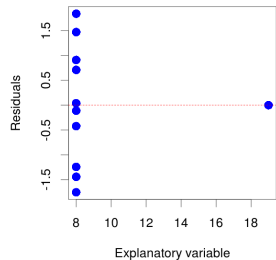
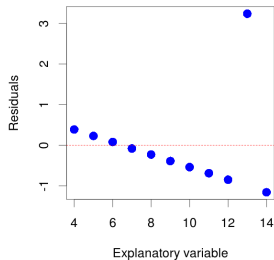
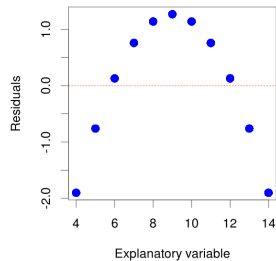
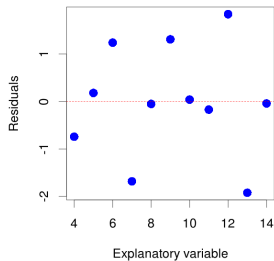
Anscombe data - Linear models

- $\bar{x} = 9$ and $\bar{y} = 7.5$ for each dataset
- $s_x^2 = 11$, $s_y^2 = 4.1$
- For all data sets, fitted regression is same:

- $\hat{Y} = 3.0 + 0.5 * X$

- All models have $R^2 = 0.67$, $\sigma^2 = 1.24$ and slope coefficients are significant at $< 1\%$ level

Checking assumptions graphically - Anscombe data residual plots



Exploratory Data Analysis (EDA)

```
1 head(mtcars)
```

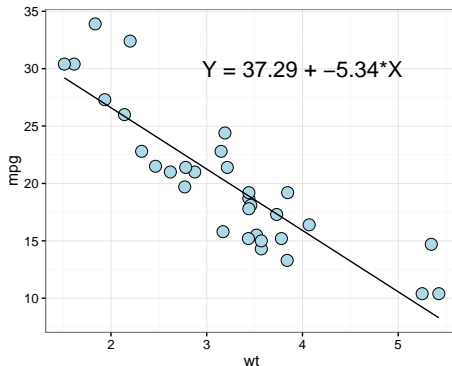
		mpg	cyl	disp	hp	drat	wt	qsec	vs	am	gear	carb
1												
2	Mazda RX4	21.0	6	160	110	3.90	2.62	16.5	0	1	4	4
3	Mazda RX4 Wag	21.0	6	160	110	3.90	2.88	17.0	0	1	4	4
4	Datsun 710	22.8	4	108	93	3.85	2.32	18.6	1	1	4	1
5	Hornet 4 Drive	21.4	6	258	110	3.08	3.21	19.4	1	0	3	1
6	Hornet Sportabout	18.7	8	360	175	3.15	3.44	17.0	0	0	3	2
7	Valiant	18.1	6	225	105	2.76	3.46	20.2	1	0	3	1

Fitting, saving and working with a linear model

```
1 # Example:
2 fit <- with(mtcars, lm(mpg ~ wt))
3 summary(fit)
4 confint(fit)
```

```
1 Call:
2 lm(formula = mpg ~ wt)
3
4 Residuals:
5     Min       1Q   Median       3Q      Max
6  -4.543  -2.365  -0.125   1.410   6.873
7
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept)   37.285     1.878   19.86 < 2e-16 ***
11 wt           -5.344     0.559   -9.56 1.3e-10 ***
12 ---
13 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
14
15 Residual standard error: 3.05 on 30 degrees of freedom
16 Multiple R-squared:  0.753,    Adjusted R-squared:  0.745
17 F-statistic: 91.4 on 1 and 30 DF,  p-value: 1.29e-10
18
19      2.5 % 97.5 %
20 (Intercept) 33.45  41.1
21 wt         -6.49  -4.2
```

- Results of linear saved as:
 - `fit.lm <- lm(Y ~ X1 + X2, mydata)`
- `fit.lm` is a linear model **object**
- Use **extractor** functions to view
- e.g., 'summary', 'coef', 'residuals', and 'fitted'



Estimation and inference from a linear model depend on several assumptions We check these assumptions using *regression diagnostics*; Assume. . .

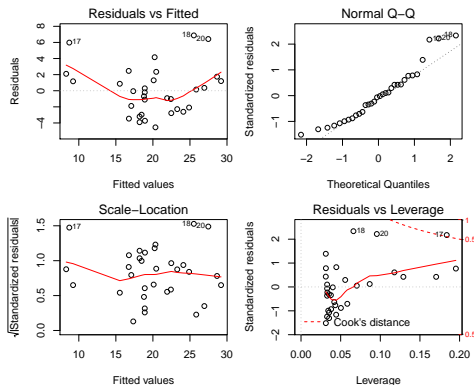
- errors are independent
- errors have constant variance
- errors are normally distributed
- all observations "fit" the model and none have large influence on the model

Violations of these assumptions can invalidate our model

Model diagnostics

- Residuals vs Fitted (*check constant variance assumption*)
- Normal Q-Q (*check normality assumption*)
- Scale-Location (*check constant variance assumption*)
- Residuals vs Leverage (*check for influential observations*)

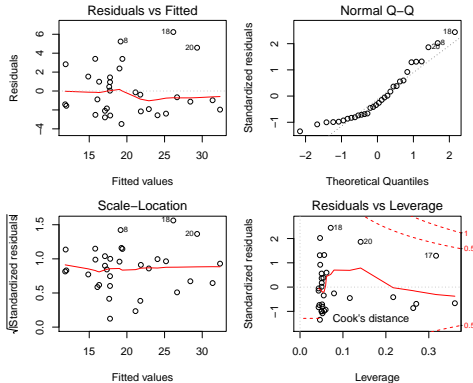
```
1 # Call the plot function on model object  
2 plot(fit)
```



Updating linear models - quadratic term

```
1 # Example:
2 fit.poly2 <- update(fit, . ~ poly(wt, degree=2))
3 summary(fit.poly2)
4 confint(fit.poly2)
```

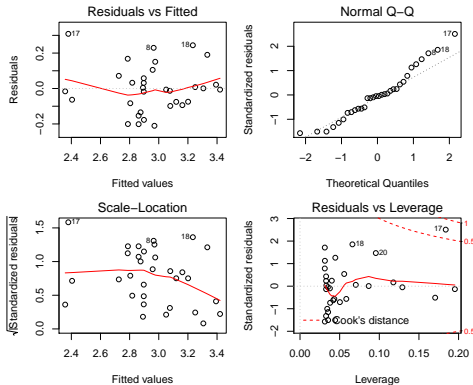
```
1 Call:
2 lm(formula = mpg ~ poly(wt, degree = 2))
3
4 Residuals:
5     Min       1Q   Median       3Q      Max
6  -3.483  -1.997  -0.773   1.462   6.238
7
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept)    20.091     0.469   42.88 < 2e-16 ***
11 poly(wt, degree = 2)1  -29.116     2.651  -10.98 7.5e-12 ***
12 poly(wt, degree = 2)2   8.636     2.651   3.26 0.0029 **
13 ---
14 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15
16 Residual standard error: 2.65 on 29 degrees of freedom
17 Multiple R-squared:  0.819,    Adjusted R-squared:  0.807
18 F-statistic: 65.6 on 2 and 29 DF,  p-value: 1.71e-11
19
20
21             2.5 % 97.5 %
22 (Intercept)    19.13    21.0
23 poly(wt, degree = 2)1  -34.54  -23.7
24 poly(wt, degree = 2)2   3.21   14.1
```



Updating linear models - transformation

```
1 # Example:
2 fit.log <- update(fit, log(mpg) ~ .)
3 summary(fit.log)
4 confint(fit.log)
```

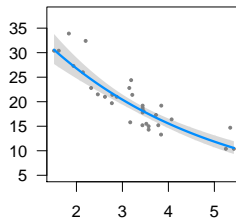
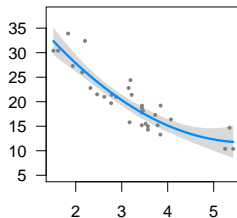
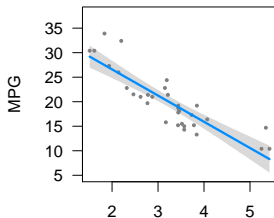
```
1
2 Call:
3 lm(formula = log(mpg) ~ wt)
4
5 Residuals:
6      Min       1Q   Median       3Q      Max
7  -0.21035 -0.08593 -0.00614  0.06133  0.30862
8
9 Coefficients:
10      Estimate Std. Error t value Pr(>|t|)
11 (Intercept)   3.832     0.084   45.6 < 2e-16 ***
12 wt           -0.272     0.025  -10.9 6.3e-12 ***
13 ---
14 Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
15
16 Residual standard error: 0.136 on 30 degrees of freedom
17 Multiple R-squared:  0.798,    Adjusted R-squared:  0.791
18 F-statistic: 118 on 1 and 30 DF, p-value: 6.31e-12
19
20      2.5 %    97.5 %
21 (Intercept)  3.660  4.003
22 wt          -0.323 -0.221
```



Visualising linear models - visreg

- Comparison of fitted models

```
1 library(visreg)
2 visreg(fit, ylim = c(5, 38), ylab = "MPG", xlab = "")
3 visreg(fit.poly2, ylim = c(5, 38), ylab = "", xlab = "WGT")
4 visreg(fit.log, "wt", trans = exp, ylim = c(5, 38), partial = TRUE, rug = FALSE, ylab = "", xlab = "")
```

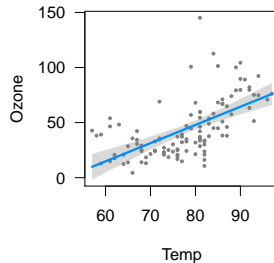
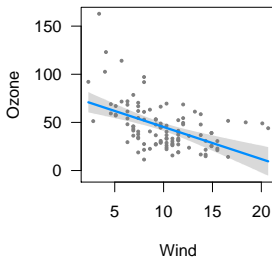
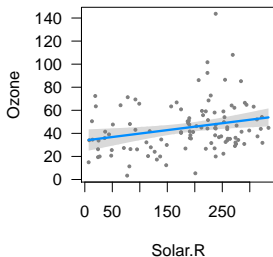


WGT

Visualising linear models - visreg

- Particularly useful when you have multiple predictors, combinations of continuous and categorical predictors, interactions

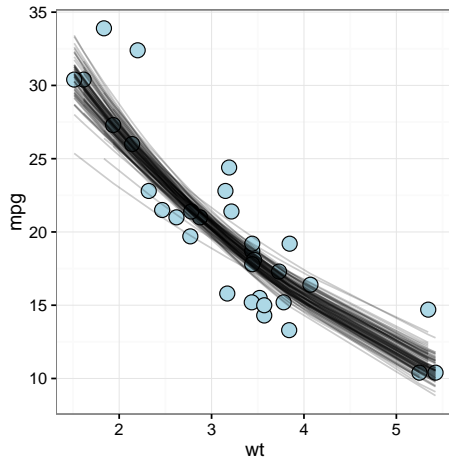
```
1 fitShow <- lm(Ozone ~ Solar.R + Wind + Temp, data = airquality)
2 visreg(fitShow)
```



The extra mile - augment and bootstrap

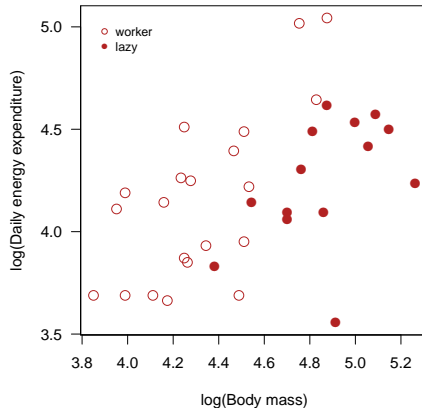
```
1 # Add variables to data frame
2 library(broom)
3 aug.mtcars <- mtcars %>% do(augment(lm(mpg ~ wt, .)))
4 head(aug.mtcars)
5 # Bootstrap model predictions and plot
6 boot.aug <- mtcars %>% bootstrap(100) %>% do(augment(lm(log(mpg) ~ wt, .)
7   ↪ ))
8 ggplot(mtcars, aes(wt, mpg)) + geom_point(size=4, shape=21, bg="lightblue"
9   ↪ ") + theme_bw() +
10   geom_line(data = boot.aug, aes(x=wt, y=exp(.fitted), group=replicate),
11     ↪ alpha=.2)
```

	.rownames	mpg	wt	.fitted	.se.fit	.resid
1	Mazda RX4	21.0	2.62	23.3	0.634	-2.283
2	Mazda RX4 Wag	21.0	2.88	21.9	0.571	-0.920
3	Datsun 710	22.8	2.32	24.9	0.736	-2.086
4	Hornet 4 Drive	21.4	3.21	20.1	0.538	1.297
5	Hornet Sportabout	18.7	3.44	18.9	0.553	-0.200
6	Valiant	18.1	3.46	18.8	0.555	-0.693



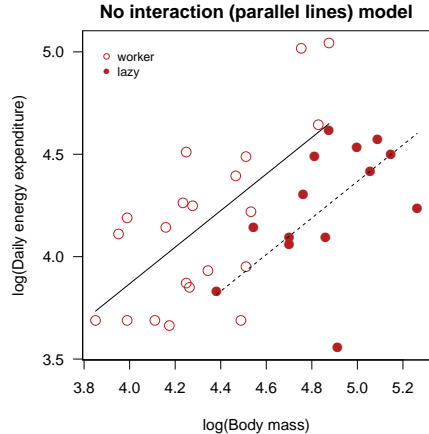
Continuous and categorical explanatory variables

- Mole rats have distinct social castes
- Single queen and small number males reproducing, remaining are workers
- Damaraland mole rats have two worker castes (worker and lazy)
- Do the groups have physiological differences?



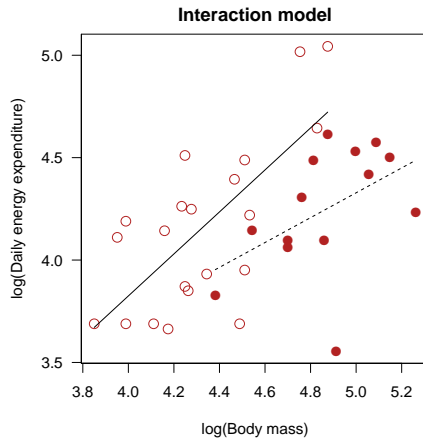
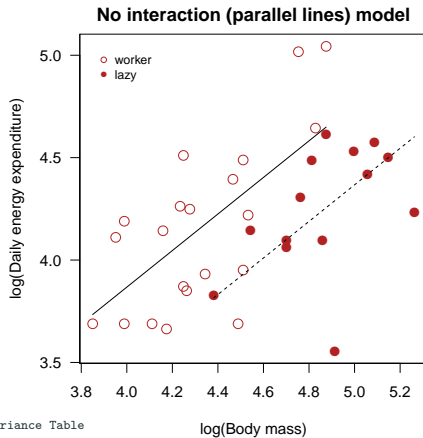
Energy-body mass relationships for mole rats

- Energy expenditure appears to vary with body mass
- But lazy workers are heavier than frequent workers
- How different is mean daily energy expenditure between caste groups when adjusted for differences in body mass?



Energy-body mass relationships for mole rats

- Need to test for interaction - non-parallel slopes



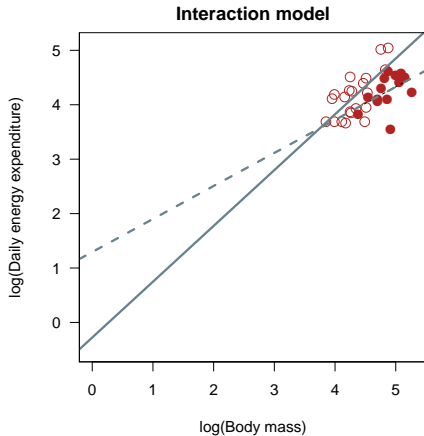
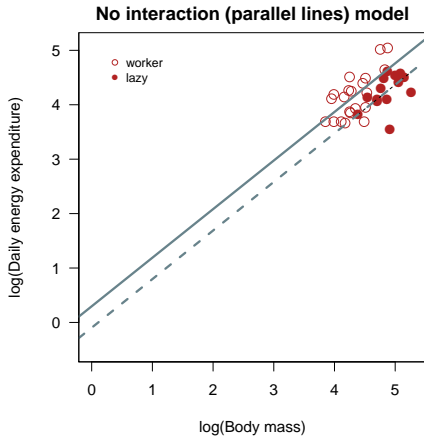
Analysis of Variance Table

Response: lnEnergy

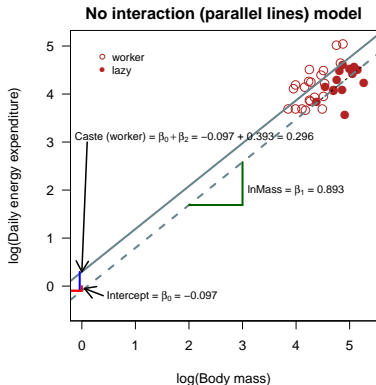
	Df	Sum Sq	Mean Sq	F value	Pr(>F)
lnMass	1	1.31	1.31	14.9	5e-04 ***
caste	1	0.64	0.64	7.3	0.01 *
lnMass:caste	1	0.09	0.09	1.0	0.32
Residuals	31	2.72	0.09		

Visualising the parameter space

- Extend plotting limits to 'see' intercepts



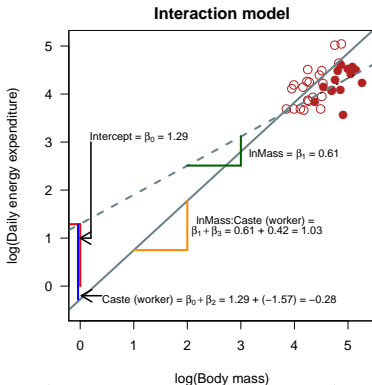
Interpretation of parameters



```
lm(lnEnergy ~ lnMass + caste)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	-0.097	0.942	-0.1	0.92
lnMass	0.893	0.193	4.6	6e-05 ***
casteworker	0.393	0.146	2.7	0.01 *

Residual standard error: 0.3 on 32 degrees of freedom
 Multiple R-squared: 0.41, Adjusted R-squared: 0.37
 F-statistic: 11 on 2 and 32 DF, p-value: 0.00022



```
lm(lnEnergy ~ lnMass + caste + lnMass*caste)
```

	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	1.29	1.67	0.8	0.44
lnMass	0.61	0.34	1.8	0.09 .
casteworker	-1.57	1.95	-0.8	0.43
lnMass:casteworker	0.42	0.41	1.0	0.32

Residual standard error: 0.3 on 31 degrees of freedom
 Multiple R-squared: 0.43, Adjusted R-squared: 0.37
 F-statistic: 7.7 on 3 and 31 DF, p-value: 0.00054

Generalised Linear Models (GLMs) can ...

- Describe relationship between mean of response and linear predictor using a *link function*
- Describe relationship between variance and mean using a *variance function*

- Sometimes our response is binary: yes/no, success/failure, diseased/non-diseased, etc. Usually coded as 1 and 0
- We use logistic regression to model the **log odds** of binary response
- Let p be the probability that response is 1. Then linear model is:

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x_1 + \cdots + \beta_k x_k$$

Converting log odds to probability

- We usually transform log odds to probability when making predictions

- Say $\log(\frac{p}{1-p}) = x$. We transform to probability as follows:

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

- R can easily do this for us (`plogis`)

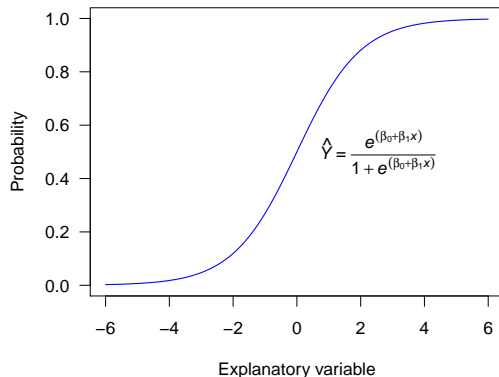
Doing Logistic Regression in R

- Logistic regression is a type of *generalized linear model* (glm)
- We fit GLM models in R using the `glm()` function. It works like the `lm()` function except we specify which glm to fit using the `family` argument
- Logistic regression requires `family=binomial`
 - Default link is `logit`, but should be specified for transparency
 - Alternatives are `probit` and `cloglog`

```
1 mod <- glm(y ~ x, data=yourdata, family=binomial)
2 summary(mod)
3 fitted(mod) # gives predicted probabilities
```

Generalised linear model for binary data - inverse link

- The logistic function has the advantage that it:
 - Yields probabilities between 0 and 1
 - Can incorporate relationships with explanatory variables through the β parameters



Generalised linear model for binary data - link function

- So our model for the probability that $Y = 1$ is the logistic function:

- $$\hat{Y} = \frac{e^{\beta_0 + \beta_1 * X_1 + \dots + \beta_p * X_p}}{1 + e^{\beta_0 + \beta_1 * X_1 + \dots + \beta_p * X_p}}$$

- But, the logistic is a nonlinear function, so we cannot estimate the β parameters using a (generalised) linear model

- Remember that we have used the inverse link $g^{-1}(\hat{Y})$

- So we have to specify this in terms of the *link function* itself, not its inverse

- This is the *logit link* (*logit* is the *log* of the *odds*)

- $$g(\hat{Y}) = \text{logit}(\hat{Y}) = \log\left(\frac{\hat{Y}}{1 - \hat{Y}}\right)$$

Example: Binomial GLM horseshoe crabs

- Nesting horseshoe crabs; female with male in nest
- Do female characteristics dictate satellite males?
- H_A : Big females have more satellites



Example: Binomial GLM horseshoe crabs

- Binary response; female with satellite male or not
- Carapace width is measure of female fitness

```
1 satt <- read.csv("../data-raw/satellites.csv", header = TRUE)
2 satt$satPresent <- (satt$nsatellites > 0)*1
3 summary(satt)
4 satt.binom <- glm(satPresent ~ width.cm, data = satt, family = binomial("
  ↳ logit"))
```

```
1      color  spine width.cm nsatellites mass.kg satPresent
2 1    medium both.bad    28.3         8    3.05         1
3 2 dark-medium both.bad    22.5         0    1.55         0
4 3 light-medium   good    26.0         9    2.30         1
5 4 dark-medium both.bad    24.8         0    2.10         0
6 5 dark-medium both.bad    26.0         4    2.60         1
7 6    medium both.bad    23.8         0    2.10         0
```

```
Call:
glm(formula = satPresent ~ width.cm, family = binomial("logit"),
    data = satt)
```

Deviance Residuals:

Min	1Q	Median	3Q	Max
-2.028	-1.046	0.548	0.907	1.694

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)
(Intercept)	-12.351	2.629	-4.70	2.6e-06 ***
width.cm	0.497	0.102	4.89	1.0e-06 ***

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

(Dispersion parameter for binomial family taken to be 1)

Null deviance: 225.76 on 172 degrees of freedom
Residual deviance: 194.45 on 171 degrees of freedom
AIC: 198.5

Number of Fisher Scoring iterations: 4

	2.5 %	97.5 %
(Intercept)	-17.810	-7.457
width.cm	0.308	0.709

Single term deletions

Model:

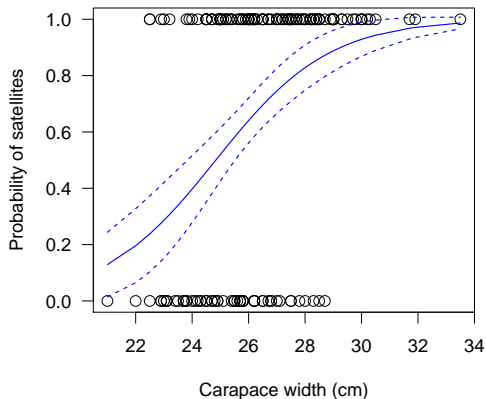
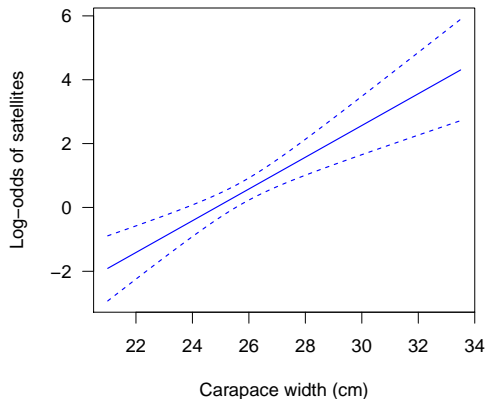
satPresent ~ width.cm

	Df	Deviance	AIC	LRT	Pr(>Chi)
<none>		194.4	198.4		
width.cm	1	225.8	227.8	31.31	2.2e-08 ***

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

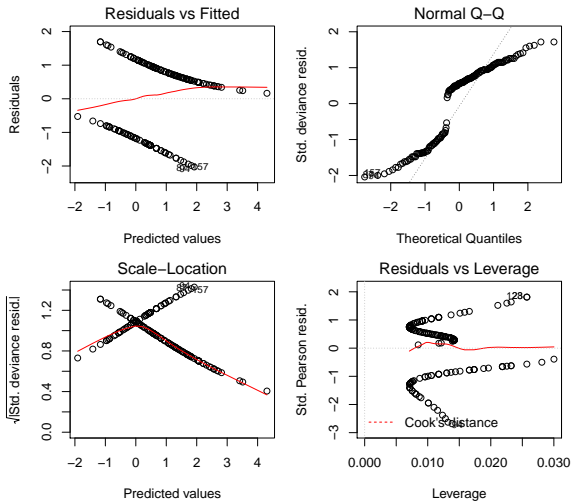
Example: Binomial GLM horseshoe crabs

- Binary response; female with satellite male or not
- Carapace width is measure of female fitness



Example: Binomial GLM horseshoe crabs

- Standard residual diagnostic plots are not informative



Challenge: Survival of passengers on the Titanic

- Sinking of the RMS Titanic is one of most infamous shipwrecks in history
- One of reasons that shipwreck led to such loss of life was that there were not enough lifeboats for passengers and crew
- Although there was some element of luck involved in surviving the sinking, some groups of people were more likely to survive than others, such as women, children, and the upper-class

Challenge question

What sorts of people were likely to survive?

This is actually part of a [kaggle competition](#)

Challenge: Survival of passengers on the Titanic

```
1 titanic <- read.csv("data-raw/titanic_long.csv")
2 head(titanic)
```

```
1   class age sex survived
2 1 first adult male      1
3 2 first adult male      1
4 3 first adult male      1
5 4 first adult male      1
6 5 first adult male      1
7 6 first adult male      1
```

