

# Chapter 10: Multiple linear regression modelling

STATS 201/8

University of Auckland

# Learning Outcomes

In this chapter you will learn about:

- Models with several explanatory variables
- Exploring pairwise relationships between all variables
- Multiple linear regression and the problem of multi-collinearity
- Fixing multi-collinearity
- Relevant R-code.

## **Section 10.1**

### **Example: Modelling birth weights using several explanatory variables**

# Multiple explanatory variables

We have learned how to model the effects of numeric and/or factor explanatory variables using linear models.

More generally, we can (in principle) fit as many explanatory variables as we like. However, we shall see that this is not always a good idea.

Caution needs to be applied.

By way of example, let us examine which variables might explain the birth weight of babies.

## Example: Birth weight of babies

<code>bwt</code>	birth weight in ounces (=28.35gm)
<code>gestation</code>	length of pregnancy in days
<code>not.first.born</code>	0=first born, 1=not first-born
<code>age</code>	mother's age in years
<code>height</code>	mother's height in inches
<code>weight</code>	mother's pre-pregnancy weight in pounds
<code>smoke</code>	smoking status of mother 0=not, 1=smoker.

The response variable is the baby's birth weight (`bwt`).

This dataset<sup>1</sup> was obtained from

<http://www.stat.berkeley.edu/users/statlabs/labs.html>.

The dataset has 1174 observations.

---

<sup>1</sup>It accompanies the excellent text Stat Labs: Mathematical Statistics through Applications Springer-Verlag (2001) by Deborah Nolan and Terry Speed.

## **Section 10.2**

### **Exploring relationships between the variables**

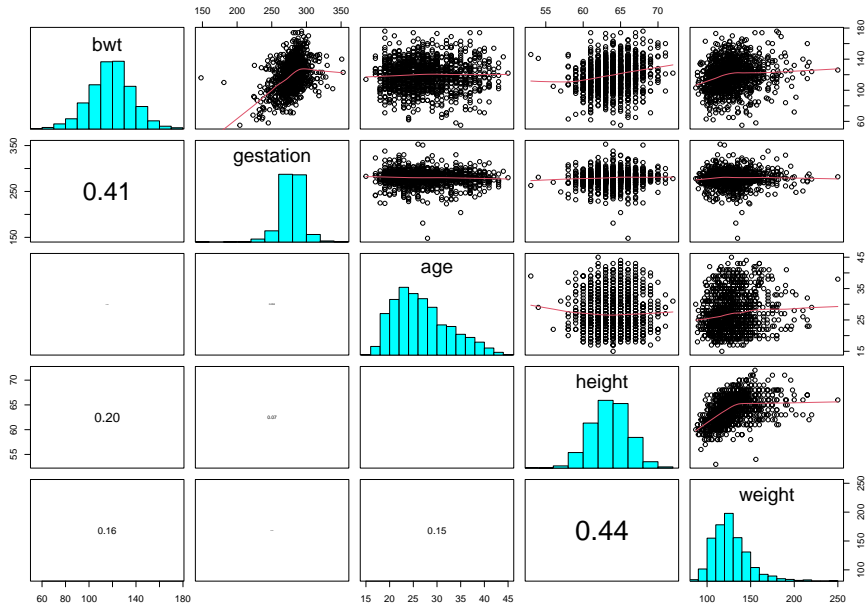
# Birth weight of babies

Let us first inspect the relationships between the numerical explanatory variables and the response variable. The numeric explanatories are *gestation*, *age*, *height* and *weight*.

The five variables are in columns 1,2,4,5 and 6 in the data frame *Babies.df*.

```
> ## Invoke the s20x library
> library(s20x)
> ## Importing data into R
> Babies.df = read.table("Data/babies_data.txt", header=T)
> ## Create the pairs plot of the five numeric variables
> pairs20x(Babies.df[,c(1,2,4,5,6)])
```

# Birth weight of babies. . .

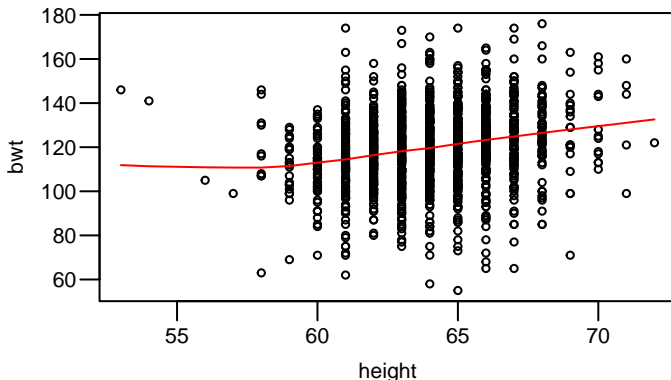




## Birth weight of babies. . .

`pairs20x` gives a histogram of each variable in the diagonal cells. Above the diagonal, in the  $(i,j)$  cell ( $i < j$ ) it gives scatter plots of variable  $i$  (y-axis) against variable  $j$  (x-axis). To illustrate, variable 1 is `bwt` and variable 4 is the mother's height (`height`). The scatter plot in cell (1,4) is

```
> plot(bwt ~ height, data = Babies.df)
> lines(lowess(Babies.df$height, Babies.df$bwt))
```



## Birth weight of babies. . .

The correlation coefficient between `height` and `bwt` is in cell (4,1). It is **0.20**, indicating that a straight-line relationship is, at best, weak.

This correlation coefficient can only measure the strength of a straight line relationship between `x` (`height`) and a `y` (`bwt`). It can be useful but can, on occasion, be misleading. In other words, look at the scatter plot and use it only if the relationship looks like a straight line.

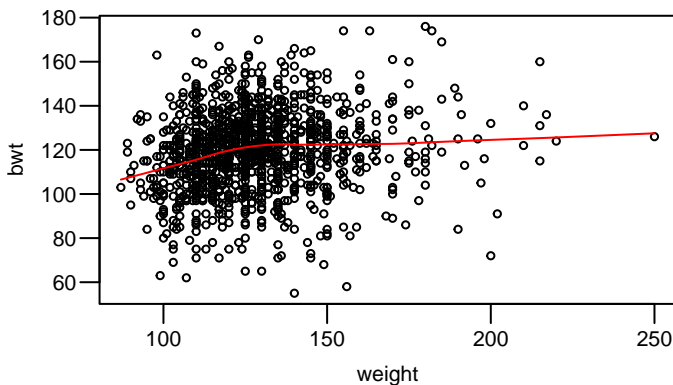
**Note:** In a simple linear regression of `y` on `x`, the resulting  $R^2$  value is the square of the sample correlation coefficient. To illustrate:

```
> summary(lm(bwt ~ height, data = Babies.df))$r.squared
[1] 0.04149539
> cor(Babies.df$bwt, Babies.df$height)^2
[1] 0.04149539
```

## Birth weight of babies. . .

Looking at the pairs plot again, we also see a somewhat weak relationship between `bwt` and mother's `weight`.

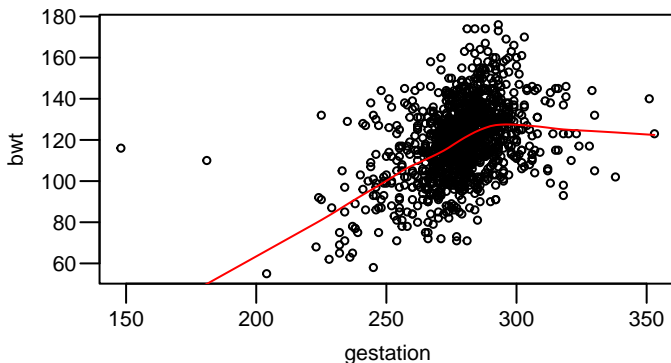
```
> plot(bwt ~ weight, data = Babies.df)
> lines(lowess(Babies.df$weight, Babies.df$bwt))
```



## Birth weight of babies. . .

There is a stronger relationship between the **gestation** time for the babies and its **bwt** which is not surprising, as the longer the child is in the mother's womb the longer the child has had time to have nutrition and grow. But, this relationship distinctly flattens out beyond a certain gestational age – some people call this a “hockey stick” curve.

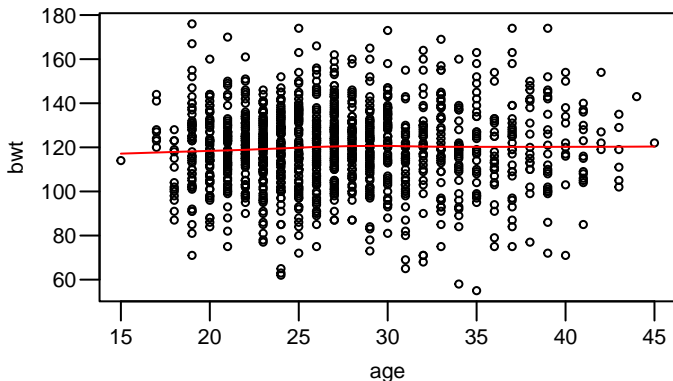
```
> plot(bwt ~ gestation, data = Babies.df)
> lines(lowess(Babies.df$gestation, Babies.df$bwt))
```



## Birth weight of babies. . .

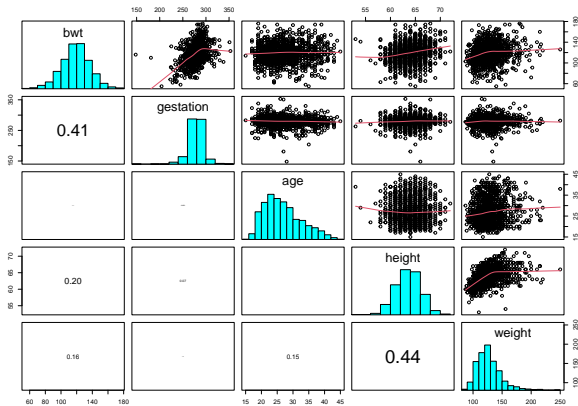
There does not seem to be any relationship between a mother's age and her child's **bwt**.

```
> plot(bwt ~ age, data = Babies.df)
> lines(lowess(Babies.df$age, Babies.df$bwt))
```



## Birth weight of babies. . .

**Note:** There seem to be some outlying data points in these plots. There does not appear to be much of a relationship between the  $x$  variables, except between **height** and **weight**.

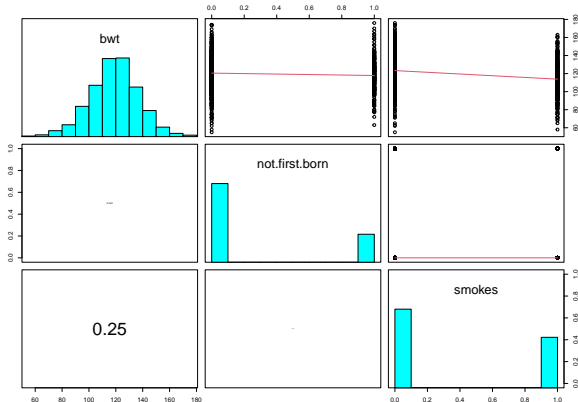


## Birth weight of babies. . .

Let us look at the categorical (factor) explanatory variables against the baby's birth weight **bwt**.

The categorical variables are **not.first.born** and **smoke**, in columns 3 and 7 of the data frame **Babies.df**.

```
> pairs20x(Babies.df[,c(1,3,7)])
```



## Birth weight of babies. . .

We see a slight decrease in babies `bwt` if the mother smokes. This increases the chance of a mother having a low birth weight baby if she smokes – perhaps another reason to avoid tobacco!

The variable `not.first.born` does not appear to have too much of an effect. This is perhaps not a surprise given that this variable may not be as important as it once was as family size has decreased markedly in the developed world (this is US data) and prenatal care has improved.

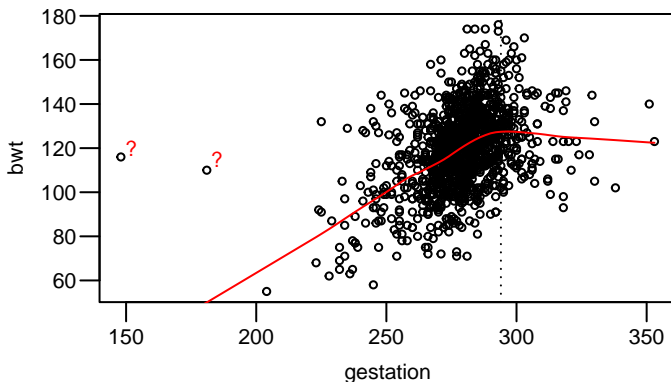
We will now begin our linear modelling of these data...



## Birth weight of babies. . .

Let us start with an understanding of **gestation** to explain **bwt** since it is the strongest relationship. The atypical data points have been marked with question marks. We will add other explanatory variables later.

```
> plot(bwt ~ gestation, data = Babies.df)
> lines(lowess(Babies.df$gestation, Babies.df$bwt), col = "red")
> text(c(152, 185), c(120, 115), "?", col = "red")
> abline(v = 294, lty = 3)
```



## Birth weight of babies. . .

Let us identify the two points denoted by the '?' symbol.

We can easily identify them in the plot as they have `gestation < 200`.

They look extremely implausible as they have typical birth-weight but have a gestational age that is extremely low for these data.

```
> id=(Babies.df$gestation<200)
> Babies.df[id,]
      bwt gestation not.first.born age height weight smokes
239 116      148           0  28     66    135      0
820 110      181           0  27     64    133      0
```

These points (observations 239 and 820) may be be unduly influential.

## Birth weight of babies. . .

The above plot has a vertical line at 294 days. The relevance of 294 days is explained in the article [“How Your Baby Grows During Pregnancy”](#).

Most babies are born before 42 weeks =  $42 \times 7 = 294$  days. It seems that beyond this point babies cease to grow and hence the ‘flattening out’ and/or decrease. In other words, it looks like the effect of gestational age depends on whether the baby is overdue or not. That is, the effect of gestational age appears to **interact** with overdue status.

It would make sense to use overdue status as an explanatory variable. Note that gestational age is a numeric variable, and being overdue (or not) is a factor variable. We already know how to fit a model using both of these as explanatory variables. Let’s create a factor variable **OD** for overdue status.

```
> Babies.df$OD = factor(as.numeric(Babies.df$gestation > 294))
> ## Check
> range(Babies.df$gestation[Babies.df$OD==0])
[1] 148 294
> range(Babies.df$gestation[Babies.df$OD==1])
[1] 295 353
```

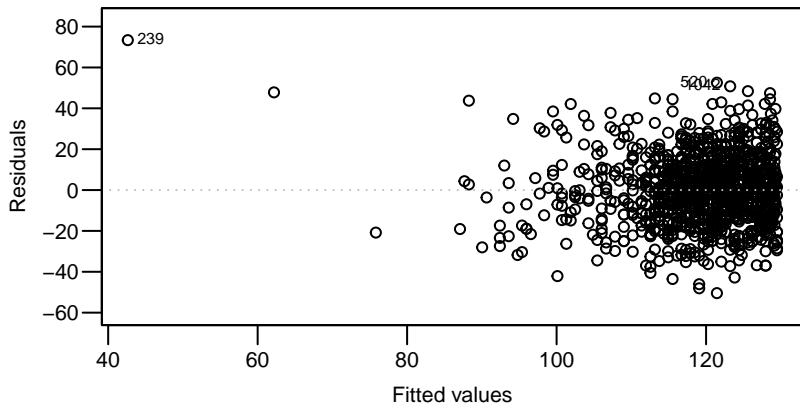
## Section 10.3

### Fitting the initial model

## Birth weight of babies. . .

Our initial model will be an interaction model using explanatory variables `gestation` and `OD`.

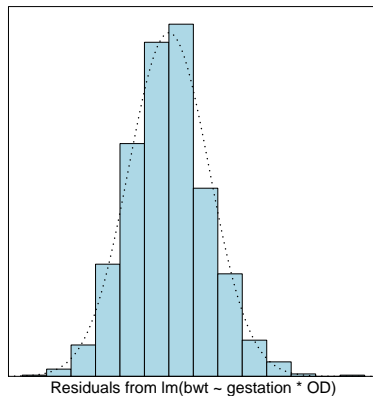
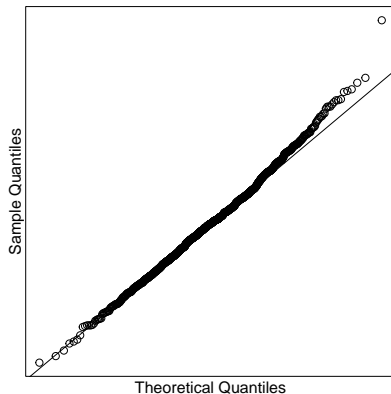
```
> bwt.fit=lm(bwt~ gestation*OD,data = Babies.df)
> plot(bwt.fit, which = 1, add.smooth = FALSE)
```



Observation 239 is a problem.

## Birth weight of babies. . .

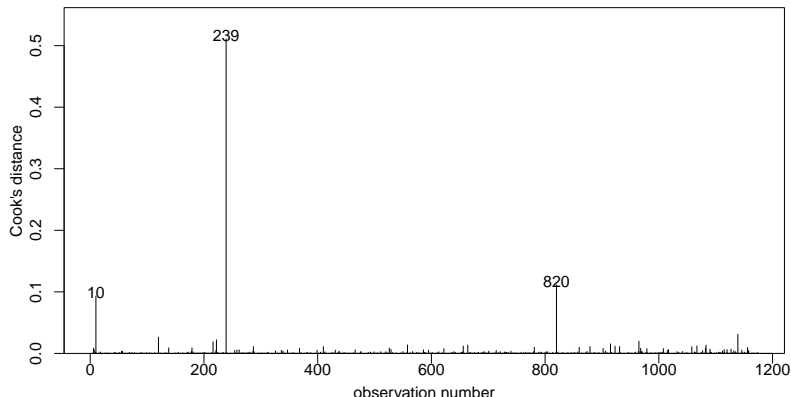
```
> normcheck(bwt.fit)
```



Other than observation 239, things look pretty good.

## Birth weight of babies. . .

```
> cooks20x(bwt.fit)
```

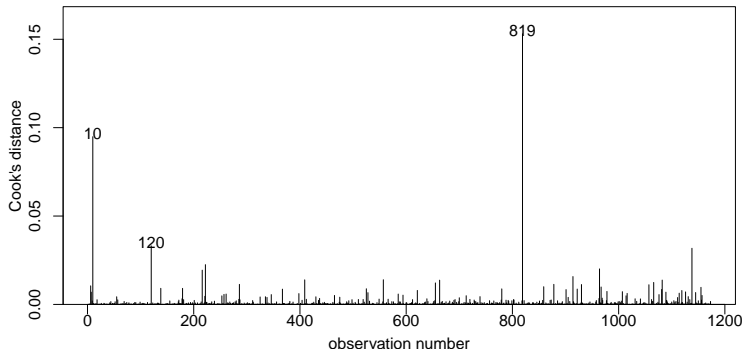


Point 239 is unduly influential. This baby has a gestational age of just 148 days, and yet has a weight typical of a full term baby. It is clearly a data-entry mistake and we will remove this data point.

## Birth weight of babies. . .

Let us refit with observation 239 removed.

```
> bwt.fit2=lm(bwt~ gestation*OD,data = Babies.df[-239,])  
> cooks20x(bwt.fit2)
```



Although observation 820<sup>2</sup> is not unduly influential, but we shall make a judgement call, and remove it.

---

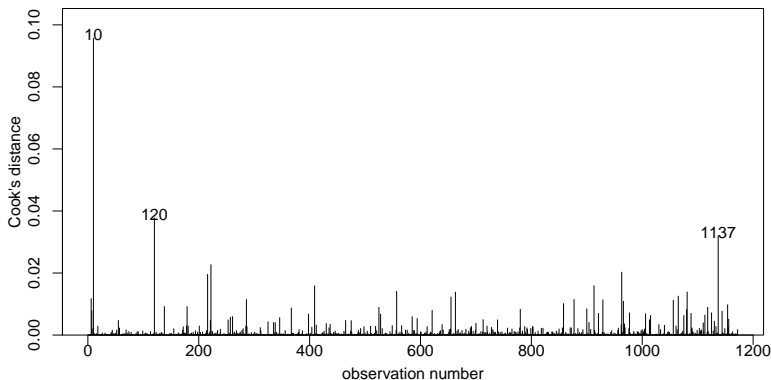
<sup>2</sup>Note that it is now identified as point 819 in this plot, but it was point 820 before we dropped point 239.



# Birth weight of babies. . .

We refit the model using the reduced data.

```
> #This time we demonstrate using the subset argument to remove points  
> bwt.fit3=lm(bwt~ gestation*OD,data = Babies.df, subset = -c(239, 820))  
> cooks20x(bwt.fit3)
```

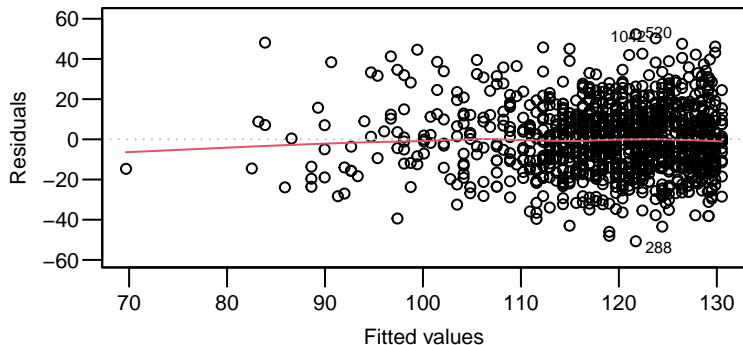


Now we have no unduly influential data points.

## Birth weight of babies. . .

Let us recheck the residuals now that we have removed these two points.

```
> plot(bwt.fit3,which=1)
```



EOV seems fine now, and the residuals seem to be centred around zero.

## Birth weight of babies. . .

Now we can trust this output let us interpret it.

```
> summary(bwt.fit3)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-68.19428	10.95985	-6.222	6.82e-10	***
gestation	0.67588	0.03963	17.054	< 2e-16	***
OD1	258.11719	41.93954	6.155	1.03e-09	***
gestation:OD1	-0.88073	0.13851	-6.358	2.92e-10	***

---

Residual standard error: 16.23 on 1168 degrees of freedom

Multiple R-squared: 0.2189, Adjusted R-squared: 0.2169

F-statistic: 109.1 on 3 and 1168 DF, p-value: < 2.2e-16

The fitted model is:

$$E[\text{bwt}] = -68.19 + 0.68 \times \text{gestation} + 258.12 \times \text{OD} \\ - 0.88 \times \text{OD} \times \text{gestation}$$

## Birth weight of babies. . .

So, for  $\text{gestation} \leq 294$  days (i.e.,  $\text{OD} = 0$ )

$$E[\text{bwt}] = -68.19 + 0.68 \times \text{gestation}$$

So, on average, babies initially grow at 0.68 oz per day until about 131 oz<sup>3</sup> at week 42 (i.e., day 294).

For  $\text{gestation} > 294$  days ( $\text{OD} = 1$ )

$$\begin{aligned} E[\text{bwt}] &= -68.19 + 258.12 + (0.68 - 0.88) \times \text{gestation} \\ &= 189.92 - 0.2 \times \text{gestation} \end{aligned}$$

So, on average, it is estimated that overdue babies lose about 0.2 oz per day after week 42.<sup>4</sup>

---

<sup>3</sup>  $131 \approx -68.19 + 0.68 \times 294$

<sup>4</sup> **Question:** How could we test whether this is significantly different from zero?

## Birth weight of babies. . .

**Note** that this model only explains about 22% of the variation in babies' birth weight, so it would be worth seeing if adding the other explanatory variables will help explain more.

In the `pairs20x` plot above we saw that `height` and `weight` had correlations of 0.20 and 0.16 with `bwt`.

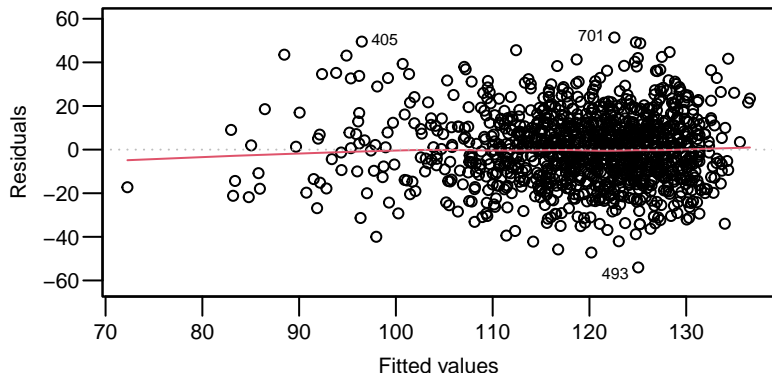
So let us see what we find when we introduce the `height` variable into the model. We will proceed with selecting variables one at a time (with reflection) – this is one of many multiple regression strategies!

**Section 10.4**  
**Multiple linear regression model:**  
**Adding more terms to the model and the peril of**  
**multi-collinearity**

## Birth weight of babies. . .

Let us add the **height** variable and see how it works out.

```
> bwt.fit4 = lm(bwt ~ gestation * OD + height, data = Babies.df,  
+ subset = -c(239,820))  
> plot(bwt.fit4, which=1)
```



All seems okay. Let us make sure that this makes sense in terms of output.

## Birth weight of babies. . .

```
> summary(bwt.fit4)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-140.41609	15.41998	-9.106	< 2e-16	***
gestation	0.65674	0.03905	16.818	< 2e-16	***
OD1	253.17362	41.21521	6.143	1.11e-09	***
height	1.21075	0.18502	6.544	8.96e-11	***
gestation:OD1	-0.86386	0.13612	-6.346	3.15e-10	***
---					

Residual standard error: 15.95 on 1167 degrees of freedom  
Multiple R-squared: 0.2465, Adjusted R-squared: 0.2439  
F-statistic: 95.45 on 4 and 1167 DF, p-value: < 2.2e-16

This seems to make sense, whereby mother's height is positively related to a baby's birth weight (on average).

**Note:** We will drop the checking of fitted vs residuals plots as it has been okay to date and it is starting to get a little tedious. We will recheck this once we get to the final model.



## Birth weight of babies. . .

Let us add `weight` to the model.

```
> bwt.fit5 = lm(bwt ~ gestation * OD + height + weight, data = Babies.df,  
+ subset = -c(239,820))  
> summary(bwt.fit5)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-132.96277	15.50163	-8.577	< 2e-16	***
gestation	0.66108	0.03889	16.998	< 2e-16	***
OD1	258.01065	41.04989	6.285	4.61e-10	***
height	0.90454	0.20460	4.421	1.07e-05	***
weight	0.08541	0.02486	3.436	0.000612	***
gestation:OD1	-0.88049	0.13558	-6.494	1.23e-10	***
---					

Residual standard error: 15.88 on 1166 degrees of freedom

Multiple R-squared: 0.2541, Adjusted R-squared: 0.2509

F-statistic: 79.43 on 5 and 1166 DF, p-value: < 2.2e-16

This makes sense. Heavier mothers can be expected to have heavier babies.

## Birth weight of babies. . .

The mother being very underweight or excessively overweight can have negative effects on their babies health, but neither **height** or **weight** directly measures this.

We will construct a new variable, body mass index **bmi**.

$$BMI = \frac{\text{mass in kg}}{\text{height in metres}^2} = \frac{\text{mass in lb}}{\text{height in inches}^2} \times 703$$

The World Health Organisation classifies BMIs in the range 18.5–25 as healthy, 25–30 as overweight, and 30+ as obese.

## Birth weight of babies. . .

Let us add `bmi` to the current model.

```
> # Create the variable BMI and add it to the model
> Babies.df$bmi = (Babies.df$weight / (Babies.df$height^2) ) * 703
> bwt.fit6 = lm(bwt ~ gestation * OD + weight + height + bmi, data = Babies.df,
+ subset = -c(239, 820))
> summary(bwt.fit6)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-217.97090	79.76837	-2.733	0.00638	**
gestation	0.66127	0.03889	17.004	< 2e-16	***
OD1	258.08808	41.04678	6.288	4.55e-10	***
weight	-0.24369	0.30395	-0.802	0.42287	
height	2.23309	1.23990	1.801	0.07196	.
bmi	1.91588	1.76352	1.086	0.27753	
gestation:OD1	-0.88083	0.13557	-6.497	1.21e-10	***
---					

Residual standard error: 15.87 on 1165 degrees of freedom

Multiple R-squared: 0.2548, Adjusted R-squared: 0.251

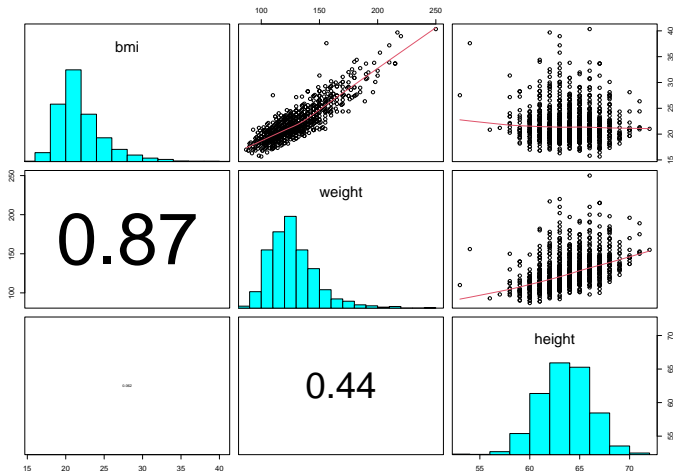
F-statistic: 66.4 on 6 and 1165 DF, p-value: < 2.2e-16

Hang on. Everything has gone weird!!! None of `weight`, `height` or `bmi` is statistically significant (at the 5% level). So what is going on?

# Birth weight of babies. . .

Let's look at these three variables to see what is happening.

```
> pairs20x(Babies.df[-c(239,820), c(9,6,5)])
```



Not surprisingly, we see that **bmi** and **weight** seem to explain each other.

## Birth weight of babies. . .

The problem is that we have a redundancy in our explanatory variables. Here, `bmi` is explained by `weight` and vice-versa. Note that adding `bmi` to the model barely changed  $R^2$  and so is telling us that it did not increase our ability to explain variability in birth weight.

In essence the statistical significance (i.e.,  $P$ -value) of an explanatory variable is measuring its contribution toward explaining variability in the response variable (in our case `bwt`) *having adjusted for any other explanatory variables in the model*.

So `bmi` explains little variability in `bwt` since `weight` has already explained most of that variability, and vice-versa.

This problem is given the name **multi-collinearity**<sup>5</sup>.

In linear algebra, we say we have linear dependence (as opposed to linear independence) in these variables.

---

<sup>5</sup>The double 'l' is not a mistake.

## Birth weight of babies...

Back to the drawing board. Let us refit this model with **bmi** and **height**, but without **weight**.

```
> bwt.fit7 = lm(bwt ~ gestation*OD + height + bmi, data = Babies.df,  
+ subset = -c(239, 820))  
> summary(bwt.fit7)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-155.30413	15.91972	-9.755	< 2e-16	***
gestation	0.66121	0.03888	17.005	< 2e-16	***
OD1	258.12344	41.04046	6.289	4.50e-10	***
height	1.25008	0.18447	6.777	1.95e-11	***
bmi	0.50673	0.14421	3.514	0.000459	***
gestation:OD1	-0.88089	0.13555	-6.499	1.20e-10	***
---					

Residual standard error: 15.87 on 1166 degrees of freedom  
Multiple R-squared: 0.2544, Adjusted R-squared: 0.2512  
F-statistic: 79.57 on 5 and 1166 DF, p-value: < 2.2e-16

Let us next investigate whether the categorical variable (**smokes**) helps to explain further variability in **bwt**.

# Birth weight of babies. . .

Let us add **smokes** to this analysis.

```
> bwt.fit8=lm(bwt~ gestation*OD+height+bmi+smokes,data = Babies.df,  
+             subset = -c(239, 820))  
> summary(bwt.fit8)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-145.23998	15.49338	-9.374	< 2e-16	***
gestation	0.63631	0.03784	16.814	< 2e-16	***
OD1	250.86546	39.83512	6.298	4.28e-10	***
height	1.28076	0.17905	7.153	1.50e-12	***
bmi	0.41557	0.14035	2.961	0.00313	**
smokes	-7.93471	0.92747	-8.555	< 2e-16	***
gestation:OD1	-0.85587	0.13157	-6.505	1.15e-10	***
---					

Residual standard error: 15.4 on 1165 degrees of freedom  
Multiple R-squared: 0.2985, Adjusted R-squared: 0.2949  
F-statistic: 82.61 on 6 and 1165 DF, p-value: < 2.2e-16

As we might have suspected, a mother smoking is associated with decreased birth weight.

# Birth weight of babies...

Let us see if `not.first.born` is useful:

```
> bwt.fit9=lm(bwt~ gestation*OD+height+bmi+smokes+not.first.born,  
+             data = Babies.df, subset = -c(239, 820))  
> summary(bwt.fit9)
```

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )	
(Intercept)	-146.36091	15.42725	-9.487	< 2e-16	***
gestation	0.64424	0.03775	17.067	< 2e-16	***
OD1	256.69266	39.69307	6.467	1.47e-10	***
height	1.29903	0.17832	7.285	5.93e-13	***
bmi	0.35512	0.14084	2.521	0.011821	*
smokes	-7.98064	0.92340	-8.643	< 2e-16	***
not.first.born	-3.50274	1.03078	-3.398	0.000701	***
gestation:OD1	-0.87488	0.13110	-6.673	3.86e-11	***

---

Residual standard error: 15.33 on 1164 degrees of freedom

Multiple R-squared: 0.3054, Adjusted R-squared: 0.3012

F-statistic: 73.1 on 7 and 1164 DF, p-value: < 2.2e-16

Hmmm, does the negative effect of `not.first.born` seem reasonable???

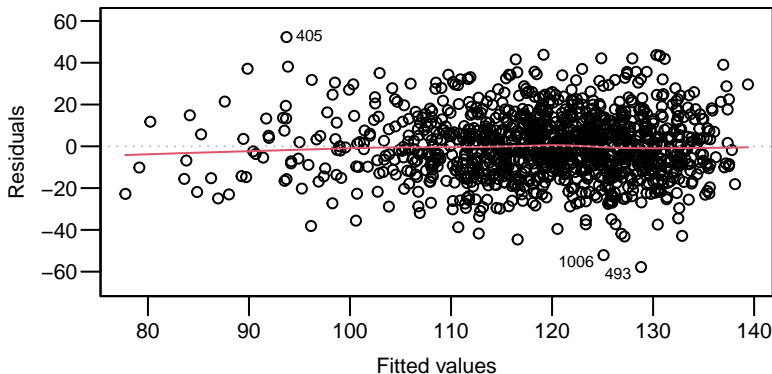


## Birth weight of babies. . .

Let us check the assumptions on this final model:

Independence should be okay, as this is (hopefully) a random sample of data from a carefully designed study.

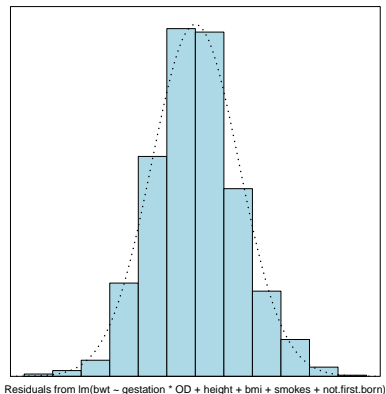
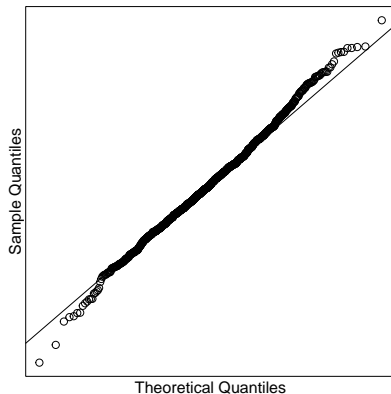
```
> plot(bwt.fit9,which=1)
```



No trend, and EOv assumption is fine.

# Birth weight of babies. . .

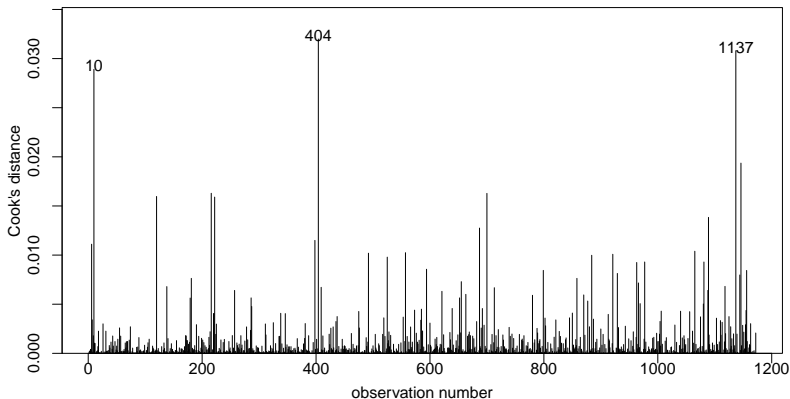
```
> normcheck(bwt.fit9)
```



Normality assumption looks fine.

## Birth weight of babies. . .

```
> cooks20x(bwt.fit9)
```



No unduly influential points.

## Birth weight of babies. . .

Let us get the CIs on this trusted output.

```
> confint(bwt.fit9)
```

	2.5 %	97.5 %
(Intercept)	-176.62923297	-116.0925891
gestation	0.57017733	0.7182955
OD1	178.81470454	334.5706219
height	0.94916219	1.6489067
bmi	0.07878781	0.6314538
smokes	-9.79235265	-6.1689248
not.first.born	-5.52512194	-1.4803515
gestation:OD1	-1.13210128	-0.6176549

See Case Study 10.1 for a detailed executive summary.

# Birth weight of babies. . .

## Closing remarks

Recall that we can fit as many explanatory variables as we like. So, did fitting all of these explanatory variables help us describe the variability of the birth weight of babies?

	What we did	Multiple $R^2$
bwt.fit3	Added gestation * OD	21.9%
bwt.fit4	Added height	24.7%
bwt.fit5	Added weight	25.4%
bwt.fit6	Added bmi	25.5%
bwt.fit7	Dropped weight	25.4%
bwt.fit8	Added smokes	29.8%
bwt.fit9	Added not.first.born	30.5%

Our final model, `bwt.fit9`, includes explanatory variables we deemed suitable and it has a Multiple  $R^2$  of 30.5%.

## Birth weight of babies. . .

### Closing remarks. . .

In situations where there are many explanatory variables, some of which may be strongly correlated, selecting the best subset for the final model can be challenging.

Model selection is a crucial component of statistical modelling and machine learning. STATS 330 (Advanced Statistical Modelling) covers this topic in more detail, using techniques such as stepwise variable selection, AIC (Akaike's information criterion), and assessment of prediction error using cross validation.

## Section 10.5

### Relevant R-code

# Most of the R-code you need for this chapter

Note that this code comes with the usual code/checks discussed in chapters 1 and 2.

Useful tools for inspecting many relationships are:

```
> ## Create the pairs plot of the five numeric variables  
> pairs20x(Babies.df[,c(1,2,4,5,6)])
```

and for the factor variables:

```
> pairs20x(Babies.df[,c(1,3,7)])
```

Then it is a dialogue with your usual outputs (e.g: `summary(bwt.fit3)` and `summary(bwt.fit4)` etc...) and to see if what you observe makes sense. Note that this requires constant vigilance.

Also note that many times several models may be selected that make sense and are acceptable.