

Lecture 5

Matrix approach to linear regression

Aim: to transition to multiple linear regression

1. Elements of multivariate normal distribution
2. Matrix approach
3. Multiple linear regression

Last week we learned

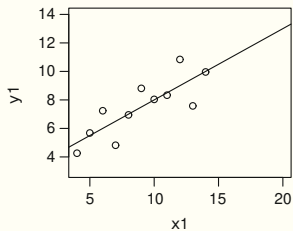
Anscombe's four data sets

Regression diagnostics

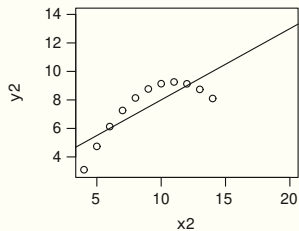
Transformations

Anscombe's four data sets

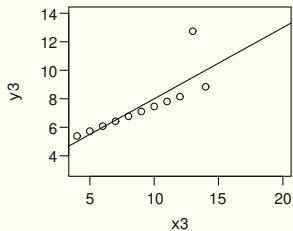
Data Set 1



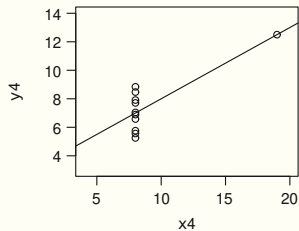
Data Set 2



Data Set 3



Data Set 4



Regression diagnostics

- Standardised residuals r_i are given by

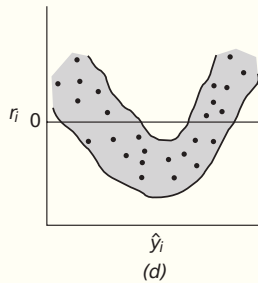
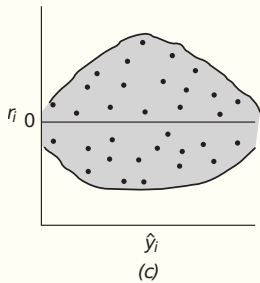
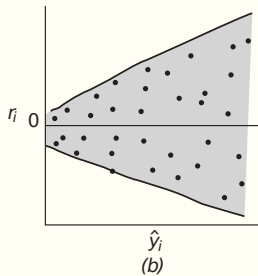
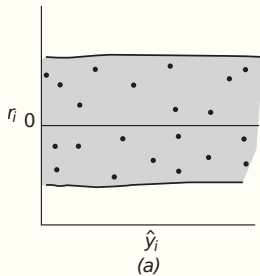
$$r_i = \frac{e_i}{\sqrt{\hat{\sigma}^2(1-h_{ii})}} \quad \text{where} \quad h_{ii} = \frac{1}{n} + \frac{(x_i - \bar{x})^2}{s_{xx}}$$

- The i th case is an outlier if
 - $|r_i| > 2$ for small and medium data sets,
 - $|r_i| > 4$ for large data sets.
- The i th case is a high leverage case if $h_{ii} > 4/n$. Moreover,
 - it is a bad leverage case if it is also an outlier,
 - it is a good leverage case otherwise.
- Cook's distance is a measure of influence given by

$$D_i = \frac{r_i^2}{2} \cdot \frac{h_{ii}}{1-h_{ii}}$$

A cut-off for D_i for the simple linear regression is $4/(n-2)$.

Constance of variance plots



Transformations

- Typical transformations are:
 - Square root transformation, i.e. $y_i \rightarrow \sqrt{y_i}$ or $x_i \rightarrow \sqrt{x_i}$
 - Log transformation, i.e. $y_i \rightarrow \log y_i$ or $x_i \rightarrow \log x_i$
 - Power transformations, i.e. $y_i \rightarrow y_i^\lambda$ or $x_i \rightarrow x_i^\lambda$

Lecture 5

Matrix approach to linear regression

Aim: to transition to multiple linear regression

1. Elements of multivariate normal distribution
2. Matrix approach
3. Multiple linear regression

Multivariate normal distribution

- Let $z_i \sim N(\mu_i, \sigma_{ii}^2)$ for $1 \leq i \leq n$ be normal random variables
- Set $\sigma_{ij}^2 = \text{Cov}(z_i, z_j)$ for $1 \leq i, j \leq n$
- The joint distribution of the z_i 's is the **multivariate normal distribution**

$$\mathbf{z} \sim N_n(\boldsymbol{\mu}, V)$$

where

$$\mathbf{z} = \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} \quad \boldsymbol{\mu} = \begin{pmatrix} \mu_1 \\ \mu_2 \\ \vdots \\ \mu_n \end{pmatrix} \quad V = \begin{pmatrix} \sigma_{11}^2 & \sigma_{12}^2 & \cdots & \sigma_{1n}^2 \\ \sigma_{21}^2 & \sigma_{22}^2 & \cdots & \sigma_{2n}^2 \\ \vdots & \vdots & \ddots & \vdots \\ \sigma_{n1}^2 & \sigma_{n2}^2 & \cdots & \sigma_{nn}^2 \end{pmatrix}$$

Expectation value

- Let \mathbf{z} be an n -dimensional column vector of random vectors, then

$$\mathbb{E}(\mathbf{z}) = \mathbb{E} \begin{pmatrix} z_1 \\ z_2 \\ \vdots \\ z_n \end{pmatrix} = \begin{pmatrix} \mathbb{E}(z_1) \\ \mathbb{E}(z_2) \\ \vdots \\ \mathbb{E}(z_n) \end{pmatrix}$$

- Let a be a scalar, \mathbf{b} – an n -dimensional column vector of constants, and A, B – matrices of constants, then

$$\mathbb{E}(a\mathbf{z} + \mathbf{b}) = a \mathbb{E}(\mathbf{z}) + \mathbf{b}$$

$$\mathbb{E}(A\mathbf{z}) = A\mathbb{E}(\mathbf{z})$$

$$\mathbb{E}(\mathbf{z}^T B) = \mathbb{E}(\mathbf{z})^T B$$

- Variance-covariance (dispersion) matrix

$$\text{Var}(\mathbf{z}) = \begin{pmatrix} \text{Var}(z_1) & \text{Cov}(z_1, z_2) & \dots & \text{Cov}(z_1, z_n) \\ \text{Cov}(z_2, z_1) & \text{Var}(z_2) & \dots & \text{Cov}(z_2, z_n) \\ \vdots & \vdots & \ddots & \vdots \\ \text{Cov}(z_n, z_1) & \text{Cov}(z_n, z_2) & \dots & \text{Var}(z_n) \end{pmatrix}$$

- The matrix $\text{Var}(\mathbf{z})$ is symmetric, $\text{Cov}(z_i, z_j) = \text{Cov}(z_j, z_i)$
- For mutually uncorrelated random variables it is a diagonal matrix, $\text{Cov}(z_i, z_j) = 0$ for all $i \neq j$
- It can be written as $\text{Var}(\mathbf{z}) = \mathbb{E}[(\mathbf{z} - \mathbb{E}(\mathbf{z}))(\mathbf{z} - \mathbb{E}(\mathbf{z}))^T]$
- For a transformed variable $\mathbf{u} = A\mathbf{z}$ we have $\text{Var}(\mathbf{u}) = A\text{Var}(\mathbf{z})A^T$

Variance-covariance matrix

- $\text{Var}(\mathbf{z}) = \mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^T]$ where $\boldsymbol{\mu} = \mathbb{E}(\mathbf{z})$. Indeed:

$$\begin{aligned}\mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^T] &= \mathbb{E}\left[\begin{pmatrix} z_1 - \mu_1 \\ z_2 - \mu_2 \\ \vdots \\ z_n - \mu_n \end{pmatrix} (z_1 - \mu_1, z_2 - \mu_2, \dots, z_n - \mu_n)\right] \\ &= \begin{pmatrix} \mathbb{E}((z_1 - \mu_1)^2) & \mathbb{E}((z_1 - \mu_1)(z_2 - \mu_2)) & \cdots & \mathbb{E}((z_1 - \mu_1)(z_n - \mu_n)) \\ \mathbb{E}((z_2 - \mu_2)(z_1 - \mu_1)) & \mathbb{E}((z_2 - \mu_2)^2) & \cdots & \mathbb{E}((z_2 - \mu_2)(z_n - \mu_n)) \\ \vdots & \vdots & \ddots & \vdots \\ \mathbb{E}((z_n - \mu_n)(z_1 - \mu_1)) & \mathbb{E}((z_n - \mu_n)(z_2 - \mu_2)) & \cdots & \mathbb{E}((z_n - \mu_n)^2) \end{pmatrix} \\ &= \text{Var}(\mathbf{z})\end{aligned}$$

- Homework: show that $\text{Var}(\mathbf{z}) = \mathbb{E}[\mathbf{z}\mathbf{z}^T] - \boldsymbol{\mu}\boldsymbol{\mu}^T$

- Let $\mathbf{u} = A\mathbf{z}$. Then $\text{Var}(\mathbf{u}) = A\text{Var}(\mathbf{z})A^T$. Indeed:

$$\begin{aligned}\text{Var}(\mathbf{u}) &= \mathbb{E}[(\mathbf{u} - \mathbb{E}(\mathbf{u}))(\mathbf{u} - \mathbb{E}(\mathbf{u}))^T] \\ &= \mathbb{E}[(A\mathbf{z} - A\boldsymbol{\mu})(A\mathbf{z} - A\boldsymbol{\mu})^T] \\ &= \mathbb{E}[A(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^T A^T] \\ &= A\mathbb{E}[(\mathbf{z} - \boldsymbol{\mu})(\mathbf{z} - \boldsymbol{\mu})^T]A^T \\ &= A\text{Var}(\mathbf{z})A^T\end{aligned}$$

Lecture 5

Matrix approach to linear regression

Aim: to transition to multiple linear regression

1. Elements of multivariate normal distribution
2. Matrix approach
3. Multiple linear regression

Matrix approach

- The SLR model has n equations

$$y_1 = \beta_0 + \beta_1 x_1 + \varepsilon_1$$

$$y_2 = \beta_0 + \beta_1 x_2 + \varepsilon_2$$

$$\vdots \qquad \qquad \vdots$$

$$y_n = \beta_0 + \beta_1 x_n + \varepsilon_n$$

- In the matrix form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$$

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad \mathbf{X} = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \boldsymbol{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \quad \boldsymbol{\varepsilon} = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

- The matrix \mathbf{X} is known as the **design matrix**.

Distribution of ε and y

- The SLR model in the matrix form

$$y = X\beta + \varepsilon$$

- Standard assumption for the error terms

$$\varepsilon_i \stackrel{\text{ind}}{\sim} N(0, \sigma^2) \implies \varepsilon \sim N_n(\mathbf{0}, \sigma^2 I_n)$$

- Distribution of y

$$y_i \stackrel{\text{ind}}{\sim} N(\beta_0 + \beta_1 x_i, \sigma^2) \implies y \sim N_n(X\beta, \sigma^2 I_n)$$

where

$$y = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} \quad X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \quad \beta = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} \quad \varepsilon = \begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix} \quad I_n = \begin{pmatrix} 1 & & & \\ & 1 & & \\ & & \ddots & \\ & & & 1 \end{pmatrix}$$

Claim. The least squares estimate of $\boldsymbol{\beta}$ is

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

Proof. First, note that

$$SS_E = \sum_i e_i e_i = \mathbf{e}^T \mathbf{e}$$

Hence in vector notation, to minimise SS_E we must solve the equation

$$\begin{pmatrix} \frac{\partial}{\partial \beta_0} \mathbf{e}^T \mathbf{e} \\ \frac{\partial}{\partial \beta_1} \mathbf{e}^T \mathbf{e} \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \implies \frac{\partial \mathbf{e}^T \mathbf{e}}{\partial \boldsymbol{\beta}} = \mathbf{0}$$

Matrix derivatives

- Let $\mathbf{z} = (z_1, \dots, z_r)^T$ and let $f(z_1, \dots, z_r)$ be a function of \mathbf{z} . Then

$$\frac{\partial f(z_1, \dots, z_r)}{\partial \mathbf{z}} = \begin{pmatrix} \frac{\partial f(z_1, \dots, z_r)}{\partial z_1} \\ \vdots \\ \frac{\partial f(z_1, \dots, z_r)}{\partial z_n} \end{pmatrix}$$

- For any vector $\mathbf{a} = (a_1, \dots, a_r)^T$ we have

$$\frac{\partial \mathbf{a}^T \mathbf{z}}{\partial \mathbf{z}} = \frac{\partial \mathbf{z}^T \mathbf{a}}{\partial \mathbf{z}} = \frac{\partial (a_1 z_1 + \dots + a_r z_r)}{\partial \mathbf{z}} = \begin{pmatrix} a_1 \\ \vdots \\ a_r \end{pmatrix} = \mathbf{a}$$

- If M is a square $r \times r$ matrix then

$$\frac{\partial \mathbf{z}^T M \mathbf{z}}{\partial \mathbf{z}} = (M + M^T) \mathbf{z}$$

Least squares estimation

- Using

$$\frac{\partial \mathbf{a}^T \mathbf{z}}{\partial \mathbf{z}} = \frac{\partial \mathbf{z}^T \mathbf{a}}{\partial \mathbf{z}} = \mathbf{a} \quad \frac{\partial \mathbf{z}^T M \mathbf{z}}{\partial \mathbf{z}} = (M + M^T) \mathbf{z}$$

we compute

$$\begin{aligned} \frac{\partial \mathbf{e}^T \mathbf{e}}{\partial \boldsymbol{\beta}} &= \frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y} - X\boldsymbol{\beta})^T (\mathbf{y} - X\boldsymbol{\beta}) \\ &= \frac{\partial}{\partial \boldsymbol{\beta}} (\mathbf{y}^T \mathbf{y} - \boldsymbol{\beta}^T X^T \mathbf{y} - \mathbf{y}^T X \boldsymbol{\beta} + \boldsymbol{\beta}^T (X^T X) \boldsymbol{\beta}) \\ &= 0 - X^T \mathbf{y} - (\mathbf{y}^T X)^T + (X^T X + (X^T X)^T) \boldsymbol{\beta} \\ &= -2X^T \mathbf{y} + 2(X^T X) \boldsymbol{\beta} \end{aligned}$$

- Next, equate the derivative to zero

$$\left. \frac{\partial \mathbf{e}^T \mathbf{e}}{\partial \boldsymbol{\beta}} \right|_{\boldsymbol{\beta}=\hat{\boldsymbol{\beta}}} = \mathbf{0} \implies (X^T X) \hat{\boldsymbol{\beta}} = X^T \mathbf{y} \implies \hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}$$

The SLR model in the matrix form

- We found that $\hat{\beta} = (X^T X)^{-1} X^T \mathbf{y}$. This must agree with our earlier results.
- The design matrix is

$$X = \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}$$

Therefore (Homework: verify the steps on this and next slide)

$$X^T \mathbf{y} = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{pmatrix} = \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix}$$

$$X^T X = \begin{pmatrix} 1 & 1 & \cdots & 1 \\ x_1 & x_2 & \cdots & x_n \end{pmatrix} \begin{pmatrix} 1 & x_1 \\ 1 & x_2 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}$$

$$\det(X^T X) = n \sum_i x_i^2 - \left(\sum_i x_i \right)^2 = n \sum_i (x_i^2 - x_i \bar{x}) = n s_{xx}$$

The SLR model in the matrix form

- The inverse of $X^T X$ is

$$(X^T X)^{-1} = \begin{pmatrix} n & \sum x_i \\ \sum x_i & \sum x_i^2 \end{pmatrix}^{-1} = \frac{1}{n s_{xx}} \begin{pmatrix} \sum x_i^2 & -\sum x_i \\ -\sum x_i & n \end{pmatrix} = \frac{1}{s_{xx}} \begin{pmatrix} \frac{1}{n} \sum x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix}$$

- Therefore

$$\begin{aligned} \hat{\beta} &= (X^T X)^{-1} X^T y \\ &= \frac{1}{s_{xx}} \begin{pmatrix} \frac{1}{n} \sum x_i^2 & -\bar{x} \\ -\bar{x} & 1 \end{pmatrix} \begin{pmatrix} \sum y_i \\ \sum x_i y_i \end{pmatrix} = \frac{1}{s_{xx}} \begin{pmatrix} \frac{1}{n} \sum x_i^2 \sum y_i - \bar{x} \sum x_i y_i \\ \sum x_i y_i - \bar{x} \sum y_i \end{pmatrix} \\ &= \frac{1}{s_{xx}} \begin{pmatrix} \frac{1}{n} \sum x_i^2 \sum y_i - \bar{x}^2 \sum y_i - \bar{x} (\sum x_i y_i - \bar{x} \sum y_i) \\ s_{xy} \end{pmatrix} \\ &= \frac{1}{s_{xx}} \begin{pmatrix} s_{xx} \bar{y} - \bar{x} s_{xy} \\ s_{xy} \end{pmatrix} = \begin{pmatrix} \bar{y} - \hat{\beta}_1 \bar{x} \\ \hat{\beta}_1 \end{pmatrix} \end{aligned}$$

which is the same result we found earlier.

Lecture 5

Matrix approach to linear regression

Aim: to transition to multiple linear regression

1. Elements of multivariate normal distribution
2. Matrix approach
3. Multiple linear regression

Multiple linear regression

- A linear regression model that has more than one predictor variable, X_1, X_2, \dots , is called a **multiple linear regression model**.
- This model generalizes the simple linear regression in two ways:
 - it allows the mean function $\mathbb{E}(Y)$ to depend on more than one predictor,
 - it can take shapes more complicated than straight lines.
- A multiple linear regression model can regress:
 - **continuous data**, e.g. weight in kg, height in cm, income in £,
 - **ordinal categorical data**, e.g.
 - level*: very low, low, medium, high, very high
 - likeness*: dislike, dislike somewhat, neutral, like somewhat, like
 - **non-ordinal categorical data**, e.g. colour or gender, blood type.

Multiple linear regression

- Suppose that a multiple linear regression model has p predictors X_1, \dots, X_p . This means that

$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p$$

- When $X_1 = x_1, \dots, X_p = x_p$ we write

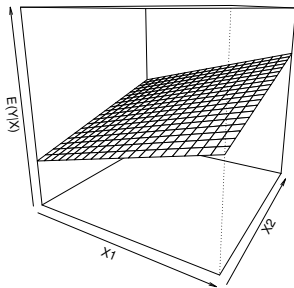
$$\mathbb{E}(Y|X = \mathbf{x}) = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p$$

- Parameters β_0, \dots, β_p are called **partial regression coefficients** – each β_i represents the expected change in the response per unit change in x_i when all x_j with $j \neq i$ are held *constant*:
 - when $p = 1$, we obtain the simple linear regression model;
 - when $p = 2$, the mean function $\mathbb{E}(Y|X)$ is a plane in 3 dimensions;
 - when $p \geq 3$, the mean function $\mathbb{E}(Y|X)$ is a hyperplane, the generalization of a p -dimensional plane in a $(p+1)$ -dimensional space.

Example: income vs education and age

- Consider the relation between the **income** and **education** of a person:
 - On an average, higher level of education provides higher income.
 - Most people have higher income when they are older than when they are young, regardless of education.
- A MLR model for income Y vs. education X_1 and age X_2

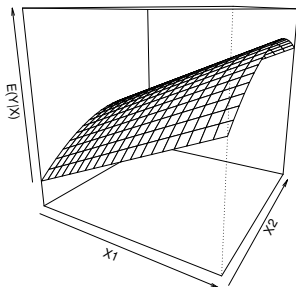
$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2$$



Example: income vs education and age

- Often the income tends to rise less rapidly in the later earning years than in early years.
- To accommodate such possibility, we might include a **quadratic term** to our model:

$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{22} X_2^2$$

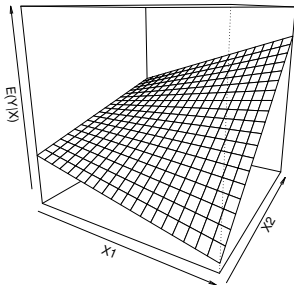


Example: a model with interactions

- A MLR model with interactions

$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_{12} X_1 X_2$$

where $X_1 X_2$ represents an interaction between X_1 and X_2 – the effect of a change in X_1 on Y depends on the level of X_2 , and vice versa



- If we let $X_3 = X_1 X_2$, we recover a MLR model with three variables.

Example: a polynomial model

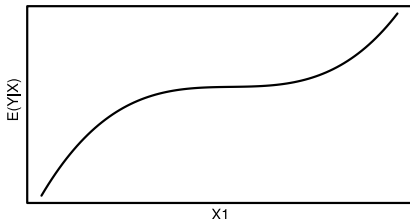
- A polynomial regression model

$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \beta_3 X_1^3$$

- The response surface is curvilinear.
- If we let $X_2 = X_1^2$, $X_3 = X_1^3$, we recover a MLR model with three variables:

$$\mathbb{E}(Y|X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3$$

- It is often advised to avoid higher order term when building the models; the reasons will be explained further in the course.



The error term

- As in the case of a simple linear regression, $\mathbb{E}(Y|X = \mathbf{x})$ does not adequately describe the data which show some randomness.
- To deal with this problem we introduce a random error $\varepsilon \sim N(0, \sigma^2)$

$$Y = \mathbb{E}(Y|X = \mathbf{x}) + \varepsilon = \beta_0 + \beta_1 x_1 + \dots + \beta_p x_p + \varepsilon$$

- MLR models in R:

```
> lm( y ~ x1 + x2 + ... + xp ) # b0 + b1 x1 + ...
```

– interaction term

```
> lm( y ~ x1*x2 ) # b0 + b1 x1 + b2 x2 + b12 x1 x2  
> lm( y ~ x1 + x2 + x1:x2 ) # the same as above
```

– polynomial model

```
> lm( y ~ x1 + I(x1^2) + I(x1^3) ) # b0 + b1 x1 + ...
```

Quiz

- Which of these regression models are linear?

(a) $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2^2 + \beta_3 x_3^3 + \varepsilon$

(b) $y = \beta_0 + \beta_1 x_1 + \beta_2^2 x_2 + \beta_3^3 x_3^3 + \varepsilon$

(c) $y = \alpha + \beta \sin(x) + \varepsilon$

(d) $y = \alpha + \log(\beta x) + \varepsilon$

(e) $y = \beta_0 + \beta_1 x_1^{-1} + \beta_2 x_2 + \beta_{123} x_1 x_2 x_3 + \varepsilon$

(f) $y = \exp(\beta_0) + \beta_1 \exp(x_1) + \exp(\beta_2 x_2) + \varepsilon$

(g) $y = X\boldsymbol{\beta} + \varepsilon$

The principle of parsimony

- Given a set of equally good explanations for a given phenomenon, then the correct explanation is the simplest explanation.
- In statistical modelling, the principle of parsimony means that:
 - models should have as few parameters as possible
 - linear models should be preferred to non-linear models
 - experiments relying on few assumptions should be preferred to those relying on many
 - models should be pared down until they are minimal adequate
 - simple explanations should be preferred to complex explanations
- In general, a variable is retained in the model only if it causes a significant increase in deviance when it is removed from the current model:
 - a model should be as simple as possible, but no simpler.

Next week

Multiple Linear Regression