

Linear statistical models The less than full rank model

flexible

Yao-ban Chan

\Rightarrow consider
each level effect of each category

The less than full rank model

In previous sections we used the linear model

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

in the knowledge (or assumption) that X , of dimension $n \times p$, is of full rank, i.e. $r(X) = p$.

This assumption is important because a full rank X implies that $X^T X$ is invertible, and therefore the normal equations

$$X^T X \mathbf{b} = X^T \mathbf{y}$$

have a unique solution.

The less than full rank model

Unfortunately, not all linear models fall into this category.

For example, consider the *one-way classification model with fixed effects*.

categorical designed variable

In this model, samples come from k distinct (sub-)populations, with different characteristics. We wish to determine the differences between these populations.

One-way classification model

For example:

- ▶ A medical researcher compares three different types of pain relievers for effectiveness in relieving arthritis;
- ▶ A botanist studies the effects of four experimental treatments used to enhance the growth of tomato plants; or
- ▶ An engineer investigates the sulfur content in the five major coal seams in a particular geographic region.

One-way classification model

Let y_{ij} be the j th sample taken from the i th population. Then the model we use is

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij}, \quad \begin{matrix} \text{mean} \\ \text{term} \end{matrix} \quad \begin{matrix} \text{population specific term} \\ \rightarrow \text{error term} \end{matrix}$$

for $i = 1, 2, \dots, k$ and $j = 1, 2, \dots, n_i$, where

- ▶ k is the number of populations/treatments;
- ▶ n_i is the number of samples from the i th population.

One-way classification model

when numerical variable:

assumption: linear relationship between design variable and response variable

1 parameter for each population

when categorical variable:

no such assumption

\Rightarrow different effect
for different category

$$\begin{bmatrix} y_{11} \\ y_{12} \\ \vdots \\ y_{21} \\ y_{22} \\ \vdots \\ y_{k,n_k} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & \dots & 0 \\ 1 & 1 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & & \vdots \\ 1 & 0 & 1 & \dots & 0 \\ 1 & 0 & 1 & \dots & 0 \\ \vdots & \vdots & \ddots & & \vdots \\ 1 & 0 & 0 & \dots & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \\ \vdots \\ \tau_k \end{bmatrix} + \begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{12} \\ \vdots \\ \varepsilon_{21} \\ \varepsilon_{22} \\ \vdots \\ \varepsilon_{k,n_k} \end{bmatrix}$$

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

The first column of X is the sum of the remaining columns, and therefore X is not of full rank.

One-way classification model

Example. Three different treatment methods for removing organic carbon from tar sand wastewater are compared: airfloatation, foam separation, and ferric-chloride coagulation. A study is conducted and the amounts of carbon removed are:

AF	FS	FCC
34.6	38.8	26.7
35.1	39.0	26.7
35.3	40.1	27.0

One-way classification model

The linear model is

$$\begin{bmatrix} 34.6 \\ 35.1 \\ 35.3 \\ 38.8 \\ 39.0 \\ 40.1 \\ 26.7 \\ 26.7 \\ 27.0 \end{bmatrix} = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \\ \tau_3 \end{bmatrix} + \begin{bmatrix} \varepsilon_{11} \\ \varepsilon_{12} \\ \varepsilon_{13} \\ \varepsilon_{21} \\ \varepsilon_{22} \\ \varepsilon_{23} \\ \varepsilon_{31} \\ \varepsilon_{32} \\ \varepsilon_{33} \end{bmatrix}$$

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

The less than full rank model

The difficulty with a less than full rank model is that $X^T X$ is singular. This means that the normal equations do not have a unique solution.

However, the problem goes deeper than that: not only can we not estimate the parameters, but the parameters themselves are not well defined.

The less than full rank model

In a one-way classification model, the response variable from population i has a mean of $\mu + \tau_i$. Thus, for our carbon removal example we might have

$$\mu + \tau_1 = 36$$

$$\mu + \tau_2 = 39$$

$$\mu + \tau_3 = 27.$$

over-parametrise

So our parameters might be $\mu = 34, \tau_1 = 2, \tau_2 = 5, \tau_3 = -7$.

However, we can also have $\mu = 30, \tau_1 = 6, \tau_2 = 9, \tau_3 = -3$.

In fact we can choose μ to be any real number, and still describe the system.

Reparametrization

One way we can tackle the less than full rank model is to convert to a full rank model. We can then use all the machinery we have developed.

Example. Consider the one-way classification model with $k = 3$.
The less than full rank model for this is

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij},$$

for $i = 1, 2, 3$, $j = 1, 2, \dots, n_i$.

However, we can write the mean of each population as

$$\mu_i = \mu + \tau_i.$$

Reparametrization

Then we can **recast the model as**

$$y_{ij} = \mu_i + \varepsilon_{ij},$$

with corresponding matrices

$$X = \begin{bmatrix} 1 & 0 & 0 \\ 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 1 & 0 \\ 0 & 1 & 0 \\ \vdots & \vdots & \vdots \\ 0 & 0 & 1 \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \mu_1 \\ \mu_2 \\ \mu_3 \end{bmatrix}.$$

$x^\tau = \begin{bmatrix} 1 & 1 & 1 & 1 & \dots \\ \dots & 1 & 1 & 1 & \dots \\ \dots & \dots & 1 & 1 & \dots \\ \dots & \dots & \dots & 1 & \dots \end{bmatrix}$

Reparametrization

The columns of X are now linearly independent, and so this is a full rank model that we can analyse. Simple matrix calculations give us

$$X^T X = \begin{bmatrix} n_1 & 0 & 0 \\ 0 & n_2 & 0 \\ 0 & 0 & n_3 \end{bmatrix}, \quad (X^T X)^{-1} = \begin{bmatrix} \frac{1}{n_1} & 0 & 0 \\ 0 & \frac{1}{n_2} & 0 \\ 0 & 0 & \frac{1}{n_3} \end{bmatrix}$$

$$X^T \mathbf{y} = \begin{bmatrix} \sum_{i=1}^{n_1} y_{1i} \\ \sum_{i=1}^{n_2} y_{2i} \\ \sum_{i=1}^{n_3} y_{3i} \end{bmatrix}, \quad \mathbf{b} = (X^T X)^{-1} X^T \mathbf{y} = \begin{bmatrix} \sum_{i=1}^{n_1} y_{1i}/n_1 \\ \sum_{i=1}^{n_2} y_{2i}/n_2 \\ \sum_{i=1}^{n_3} y_{3i}/n_3 \end{bmatrix}.$$

mean response from population i

Reparametrization

Therefore, the least squares estimates for each of the population means are the means of the samples drawn from that population:

$$\hat{\mu}_i = \bar{y}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} y_{ij}.$$

Linear functions of the parameters, of the form $\mathbf{t}^T \boldsymbol{\beta}$, are estimated using $\mathbf{t}^T \mathbf{b}$. For example, the function $\mu_1 - \mu_2$ is estimated by

$$\bar{y}_1 - \bar{y}_2 = \frac{1}{n_1} \sum_{i=1}^{n_1} y_{1i} - \frac{1}{n_2} \sum_{i=2}^{n_2} y_{2i}.$$

Reparametrization

$$SS_{\text{Res}} = SS_{\text{Total}} - SS_{\text{Reg}}$$

$$\begin{aligned} &\downarrow \\ &= (\mathbf{y} - \mathbf{X}\mathbf{b})^T (\mathbf{y} - \mathbf{X}\mathbf{b}) \\ &= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} \end{aligned}$$

The standard assumption that the errors are normally distributed with mean $\mathbf{0}$ and variance $\sigma^2 I$ is interpreted in this context to mean that all populations have a common variance σ^2 (but different means). The estimator for this variance is $\not\propto$

$$s^2 = \frac{\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}}{n - p} = \frac{\mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} \mathbf{b}}{n - p}.$$

Reparametrization

For the example,

$$\begin{aligned}
 s^2 &= \frac{1}{n-3} \left[\sum_{i=1}^3 \sum_{j=1}^{n_i} y_{ij}^2 - \left[\begin{array}{ccc} \sum_{i=1}^{n_1} y_{1i} & \sum_{i=1}^{n_2} y_{2i} & \sum_{i=1}^{n_3} y_{3i} \end{array} \right] \left[\begin{array}{c} \sum_{i=1}^{n_1} y_{1i}/n_1 \\ \sum_{i=1}^{n_2} y_{2i}/n_2 \\ \sum_{i=1}^{n_3} y_{3i}/n_3 \end{array} \right] \right] \\
 &= \frac{1}{n-3} \left[\sum_{i=1}^3 \sum_{j=1}^{n_i} y_{ij}^2 - \sum_{i=1}^3 \frac{1}{n_i} \left(\sum_{j=1}^{n_i} y_{ij} \right)^2 \right] \\
 &= \frac{1}{n-3} \sum_{i=1}^3 \left[\sum_{j=1}^{n_i} y_{ij}^2 - \frac{1}{n_i} \left(\sum_{j=1}^{n_i} y_{ij} \right)^2 \right].
 \end{aligned}$$

Reparametrization

This can be written as a 'pooled' variance

$$s^2 = \frac{(n_1 - 1)s_1^2 + (n_2 - 1)s_2^2 + (n_3 - 1)s_3^2}{(n_1 - 1) + (n_2 - 1) + (n_3 - 1)}$$

where s_i^2 are the individual population variance estimators

$$s_i^2 = \frac{1}{n_i - 1} \sum_{j=1}^{n_i} (y_{ij} - \bar{y}_{i\cdot})^2.$$

More generally, for a one-way classification model with k levels,

$$s^2 = \frac{\sum_{i=1}^k (n_i - 1)s_i^2}{\sum_{i=1}^k (n_i - 1)}.$$

Reparametrization

In general, it is always possible to re-parameterise a less than full rank model into a full rank model.

However, this is not always desirable.

For the one-way classification model, we have a nice interpretation of the (re-)parameters as the population means. But this is not always possible.

Reparametrization

Example. Consider the two-way classification model (without interaction), with two levels of each factor:

$$y_{ij} = \mu + \tau_i + \beta_j + \varepsilon_{ij}, i, j = 1, 2.$$

overall mean μ \rightarrow effect comes from τ_i
 error term ε_{ij}

The design matrix for this model is

$$X = \begin{bmatrix} 1 & \boxed{\tau_1 \quad \tau_2} & \boxed{\beta_1 \quad \beta_2} \\ \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots & \vdots \end{bmatrix}.$$

less than full rank model

\Rightarrow bnt rank = 3

\Rightarrow sum up column τ_1, τ_2

get all 1

sum up column β_1, β_2

get all 1

Reparametrization

It is obvious that the first column is the sum of the next two columns, and also the sum of the 4th and 5th columns. Thus $r(X) = 3$.

This means that we have to remove 2 parameters — which ones?
This makes interpretability much harder!

Fortunately, we do not have to re-parameterise our models: we can develop theory for the less than full rank model.

Conditional inverses

We start with more linear algebra. Here we introduce the concept of **conditional inverses**.

Definition 6.1

Let A be a $n \times p$ matrix. The $p \times n$ matrix A^c is called a *conditional inverse* for A if and only if

has the same dimension as transpose of A
+

$$AA^cA = A.$$

Conditional inverses

$$\text{then } AA^c A = A \\ AA^c A A^{-1} = A \cdot A^{-1} \\ \underline{AA^c = I.} \\ \Rightarrow A^c = A^{-1}$$

If A is square and nonsingular, then $A^{-1} = A^c$, so conditional inverses are an extension of regular inverses to non-square and singular matrices.

Example. Consider the (singular) matrices

$$A = \begin{bmatrix} 2 & 4 & 2 \\ 1 & 0 & -1 \\ 3 & 1 & -2 \end{bmatrix}, \quad A_1 = \begin{bmatrix} 0 & 1 & 0 \\ \frac{1}{4} & -\frac{1}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

Conditional inverses

We have

$$\begin{aligned}
 AA_1A &= \begin{bmatrix} 2 & 4 & 2 \\ 1 & 0 & -1 \\ 3 & 1 & -2 \end{bmatrix} \begin{bmatrix} 0 & 1 & 0 \\ \frac{1}{4} & -\frac{1}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 4 & 2 \\ 1 & 0 & -1 \\ 3 & 1 & -2 \end{bmatrix} \\
 &= \begin{bmatrix} 2 & 4 & 2 \\ 1 & 0 & -1 \\ 3 & 1 & -2 \end{bmatrix} \begin{bmatrix} 1 & 0 & -1 \\ 0 & 1 & 1 \\ 0 & 0 & 0 \end{bmatrix} \\
 &= \begin{bmatrix} 2 & 4 & 2 \\ 1 & 0 & -1 \\ 3 & 1 & -2 \end{bmatrix} = A.
 \end{aligned}$$

! conditional inverse is not a reflective relationship

\Rightarrow if A_1 is a conditional inverse of A

it doesn't mean that

A is a conditional inverse of A_1 .

∅.

Therefore A_1 is a conditional inverse for A .

Conditional inverses

how many conditional inverse can a matrix have?

1 / inf

if the matrix has an inverse, the conditional inverse should be the regular inverse (unique)

if matrix has 2 conditional inverse, then it has inf number
 \Rightarrow any linear comb

But it can also be shown that

if A has more than 1 conditional inverse, then A must be singular matrix

$$A_2 = \begin{bmatrix} 0 & 1 & 0 \\ 0 & -3 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

is also a conditional inverse for A ! So conditional inverses are not unique.

That is why we speak of (a conditional inverse) for A , not (the conditional inverse) for A .

Of course, if A is nonsingular, then the conditional inverse is uniquely the regular inverse. We can use this in the above example to show that A is singular.

$$ABA = ACA^{-1}A \Rightarrow A(2B + (r-2)C)A = 2ABA + (r-2)ACA = A.$$

Finding a conditional inverse

For a square matrix to have a regular inverse, it must satisfy some nonsingularity conditions. However, this is not the case for a conditional inverse.

$$\begin{bmatrix} M^{-1} & \\ \hline 0 & 0 \end{bmatrix}^T$$

Theorem 6.2

Let A be a $n \times p$ matrix. Then A has a conditional inverse.

Moreover, conditional inverses can be constructed as follows:

1. Find a minor M of A which is nonsingular and of dimension $r(A) \times r(A)$.
2. Replace M in A with $(M^{-1})^T$ and the other entries with zeros.
3. Transpose the resulting matrix.

minor M
 { principal minor

Finding a conditional inverse

Proof. Let's assume M is the principal (top left) minor of A . We write

$$A = \left[\begin{array}{c|c} M & A_{12} \\ \hline A_{21} & A_{22} \end{array} \right].$$

The procedure constructs a $p \times n$ matrix B which can be partitioned as

$$B = \left[\begin{array}{c|c} M^{-1} & 0 \\ \hline 0 & 0 \end{array} \right].$$

Finding a conditional inverse

Then we have

$$\begin{aligned}
 ABA &= \left[\begin{array}{c|c} M & A_{12} \\ \hline A_{21} & A_{22} \end{array} \right] \left[\begin{array}{c|c} M^{-1} & 0 \\ \hline 0 & 0 \end{array} \right] \left[\begin{array}{c|c} M & A_{12} \\ \hline A_{21} & A_{22} \end{array} \right] \\
 &= \left[\begin{array}{c|c} I & 0 \\ \hline A_{21}M^{-1} & 0 \end{array} \right] \left[\begin{array}{c|c} M & A_{12} \\ \hline A_{21} & A_{22} \end{array} \right] \\
 &= \left[\begin{array}{c|c} M & A_{12} \\ \hline A_{21} & A_{21}M^{-1}A_{12} \end{array} \right].
 \end{aligned}$$

We merely have to show that $A_{21}M^{-1}A_{12} = A_{22}$.

Finding a conditional inverse

This follows because $r(A)$ is the size of M ; we can write all other columns of A as linear combinations of the first $r(A)$ columns. In other words, there exists a matrix R such that

$$\begin{bmatrix} A_{12} \\ A_{22} \end{bmatrix} = \begin{bmatrix} M \\ A_{21} \end{bmatrix} R$$

$$A_{12} = MR$$

$$R = M^{-1}A_{12}$$

$$A_{22} = A_{21}R$$

$$= A_{21}M^{-1}A_{12}.$$

~~Not linear combination~~ \rightarrow

Finding a conditional inverse

Example. From the previous example,

$$A = \begin{bmatrix} 2 & 4 & 2 \\ 1 & 0 & -1 \\ 3 & 1 & -2 \end{bmatrix}.$$

It can be seen that $r(A) = 2$, so we take the principal 2×2 minor

$$M = \begin{bmatrix} 2 & 4 \\ 1 & 0 \end{bmatrix}.$$

Finding a conditional inverse

Then

$$(M^{-1})^T = -\frac{1}{4} \begin{bmatrix} 0 & -4 \\ -1 & 2 \end{bmatrix}^T = \begin{bmatrix} 0 & \frac{1}{4} \\ 1 & -\frac{1}{2} \end{bmatrix}$$

and

$$A^c = \begin{bmatrix} 0 & \frac{1}{4} & 0 \\ 1 & -\frac{1}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix}^T = \begin{bmatrix} 0 & 1 & 0 \\ \frac{1}{4} & -\frac{1}{2} & 0 \\ 0 & 0 & 0 \end{bmatrix}.$$

This is the conditional inverse A_1 of the earlier example, so we have seen that it works.

On the other hand, if we take the lower left 2×2 minor, following the procedure gives us A_2 . So this procedure can produce more than one conditional inverse.

Conditional inverse properties

$$n \left[\begin{array}{c} \xleftarrow{\text{P}} \\ \end{array} \right] \quad r(A) \geq r(AA^c) \geq r(AA^c A) = r(A).$$

Let A be a $n \times p$ matrix of rank r , where $n \geq p \geq r$. Then

- ▶ $r(A) = r(AA^c) = r(A^c A)$;
- ▶ $(A^c)^T = (A^T)^c$;
- ▶ $A^c A$, AA^c , $I - A^c A$ and $I - AA^c$ are idempotent;
- ▶ $A = A(A^T A)^c (A^T A)$ and $A^T = (A^T A)(A^T A)^c A^T$.

$$A^T A = A^T A (A^T A)^c (A^T A)$$

A^T may not be invertible \neq $(A^T)^{-1}$ may not exist).

Conditional inverse properties

We say that an expression involving a conditional inverse is *unique* if it is the same no matter what conditional inverse we use.

look like hat matrix

- ▶ $A(A^T A)^c A^T$ is unique, symmetric, and idempotent;
- ▶ $r(A(A^T A)^c A^T) = r$;
- ▶ $I - A(A^T A)^c A^T$ is unique, symmetric and idempotent;
- ▶ $r(I - A(A^T A)^c A^T) = n - r$.

Conditional inverse properties

Proof.

$$\begin{aligned}[A(A^T A)^c A^T]^T &= A[(A^T A)^c]^T A^T = A[(A^T A)^T]^c A^T \\ &= A(A^T A)^c A^T.\end{aligned}$$

not necessarily $= A^T A$

$$\begin{aligned}A \boxed{(A^T A)^c A^T A (A^T A)^c} A^T &= [A(A^T A)^c A^T A] (A^T A)^c A^T \\ &= A(A^T A)^c A^T.\end{aligned}$$

$r(AB) \leq \min(r(A), r(B))$.

multiple \Rightarrow rank will decrease

$$\begin{aligned}r(A(A^T A)^c A^T) &\geq r(A(A^T A)^c A^T A) \\ &= r(A) \\ &\geq r(A(A^T A)^c A^T).\end{aligned}$$

R Example

```
> library(MASS)
> A <- matrix(c(2,-6,3,1,6,4,-2,-1,0),3,3)
> det(A)

[1] 89

> A # non-singular

[,1] [,2] [,3]
[1,]    2    1   -2
[2,]   -6    6   -1
[3,]    3    4    0
```

R Example

\Rightarrow generalise inverse
 $\text{> } Ac \leftarrow \text{ginv}(A)$

$\text{> } Ac$

	[,1]	[,2]	[,3]
[1,]	0.04494382	-0.08988764	0.1235955
[2,]	-0.03370787	0.06741573	0.1573034
[3,]	-0.47191011	-0.05617978	0.2022472

$\text{> } solve(A)$

	[,1]	[,2]	[,3]
[1,]	0.04494382	-0.08988764	0.1235955
[2,]	-0.03370787	0.06741573	0.1573034
[3,]	-0.47191011	-0.05617978	0.2022472

\Rightarrow Moore - Penrose inverse

$$AA^c A = A$$

condition
+ $A^c A A^c = A^c$

condition
+ $AA^c, A^c A$ symmetric

R Example

```
> A <- matrix(c(2,-6,3,1,6,4,3,0,7),3,3)
> det(A)

[1] 0

> A # singular

[,1] [,2] [,3]
[1,]    2    1    3
[2,]   -6    6    0
[3,]    3    4    7
```

R Example

```
> Ac <- ginv(A)
> Ac

 [,1]      [,2]      [,3]
[1,] 0.025713835 -0.084240416 0.03659883
[2,] 0.009149708  0.080454330 0.04369774
[3,] 0.034863543 -0.003786086 0.08029658

> round(A %*% Ac %*% A, 5)

 [,1] [,2] [,3]
[1,]   2   1   3
[2,]  -6   6   0
[3,]   3   4   7
```

R Example

```
> Ac2 <- matrix(0,3,3)      top left minor
```

```
> Ac2[1:2,1:2] <- t(solve(A[1:2,1:2]))
```

```
> Ac2 <- t(Ac2)
```

```
> Ac2
```

```
          [,1]      [,2]  [,3]  
[1,] 0.3333333 -0.05555556 0  
[2,] 0.3333333  0.11111111 0  
[3,] 0.0000000  0.00000000 0
```

```
> A %*% Ac2 %*% A
```

```
          [,1]  [,2]  [,3]  
[1,]    2    1    3  
[2,]   -6    6    0  
[3,]    3    4    7
```

R Example

```
> Ac3 <- matrix(0,3,3)
> Ac3[2:3,1:2] <- t(solve(A[2:3,1:2]))
> Ac3 <- t(Ac3)
> Ac3

 [,1]      [,2]      [,3]
[1,]    0 -0.09523810 0.1428571
[2,]    0  0.07142857 0.1428571
[3,]    0  0.00000000 0.0000000

> A %*% Ac3 %*% A

 [,1] [,2] [,3]
[1,]   2   1   3
[2,]  -6   6   0
[3,]   3   4   7
```

R Example

```
> library(Matrix)
> rankMatrix(A)[1]
[1] 2

> rankMatrix(Ac2 %*% A)[1]
[1] 2
> round(A %*% ginv(t(A) %*% A) %*% t(A) %*% A, 5)
[,1] [,2] [,3]
[1,]    2    1    3
[2,]   -6    6    0
[3,]    3    4    7
```

$(ATA)^{-1}$

R Example

```

> A %*% ginv(t(A) %*% A) %*% t(A)
              [,1]      [,2]      [,3]
[1,]  0.16516801 -0.09938476 0.35778514
[2,] -0.09938476  0.98816848 0.04259347
[3,]  0.35778514  0.04259347 0.84666351

> AtAc2 <- matrix(0,3,3)
> AtAc2[1:2,1:2] <- solve((t(A) %*% A)[1:2,1:2])
> A %*% AtAc2 %*% t(A)
              [,1]      [,2]      [,3]
[1,]  0.16516801 -0.09938476 0.35778514
[2,] -0.09938476  0.98816848 0.04259347
[3,]  0.35778514  0.04259347 0.84666351

```

$A(A^T A)^{-1} A^T$
unique



Solving the normal equations

Let us now solve the normal equations

$$X^T X \mathbf{b} = X^T \mathbf{y}.$$

First, we must make sure that they have a solution!

Theorem 6.3

solvable , at least one set of values that satisfies each equation in the system

The system $Ax = g$ is consistent if and only if the rank of

$[A | g]$ is equal to the rank of A .

$$\text{r}([A | g]) \geq \text{r}(A)$$

Solving the normal equations

Proof. (\Leftarrow) Since $r([A | g]) = r(A)$, g must be a linear combination of the columns of A.

Therefore there exist constants x_1, x_2, \dots, x_p so that

$$x_1 \mathbf{a}_1 + x_2 \mathbf{a}_2 + \dots + x_p \mathbf{a}_p = \mathbf{g}$$

where \mathbf{a}_i is the i th column of A .

Solving the normal equations

But if we put this into matrix notation and set

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_p \end{bmatrix}, \quad \left[\begin{array}{cccc|c} 1 & 1 & \cdots & 1 & x_1 \\ a_1 & a_2 & \cdots & a_p & x_2 \\ 1 & 1 & \cdots & 1 & \vdots \\ x_p \end{array} \right]$$

then this is exactly the system $A\mathbf{x} = \mathbf{g}$.

Therefore the system is consistent.

Solving the normal equations

Theorem 6.4

In the general linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, the normal equations

$$\mathbf{X}^T \mathbf{X} \mathbf{b} = \mathbf{X}^T \mathbf{y}$$

are consistent.

Solving the normal equations

Proof. It is obvious that $r(X^T X) \leq r([X^T X \mid X^T \mathbf{y}])$, as adding a column cannot decrease the number of linearly independent columns.

However,

$$\begin{aligned}
 r([X^T X \mid X^T \mathbf{y}]) &= r(X^T [X \mid \mathbf{y}]) \\
 &\stackrel{\text{multiply}}{\leq} r(X^T) \\
 &= r(X^T X).
 \end{aligned}
 \quad \text{Handwritten note: } r(x^T) = r(x) = r(x^T x)$$

Theorem 6.3 now shows that the normal equations are consistent.

Solving the normal equations

if A non-singular
 $\Rightarrow x = A^{-1}g$

if A is singular
 $x = A^c g$

Now that we know the normal equations always have a solution,
 how can we find one?

Theorem 6.5

verify $A^c g$ is a solution

Let $Ax = g$ be a consistent system. Then $A^c g$ is a solution to the system, where A^c is any conditional inverse for A .

\Rightarrow since $Ax = g$ is a consistent system \Rightarrow there is a solution x^*

Proof. Since $Ax^* = g$ for some x^* ,

$$AA^c g = AA^c Ax^* = Ax^* = g.$$

Therefore, $A^c g$ solves the system.

Solving the normal equations

From this theorem, we see that

$$\mathbf{b} = (X^T X)^c X^T \mathbf{y}$$

solves the normal equations, for any conditional inverse.

However, in the less than full rank model, different conditional inverses may result in different solutions.

Solving the normal equations

Example. Suppose we have a one-way classification model with two classes and one sample from each class. The design matrix is

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}.$$

Supposing that $\mathbf{y}^T = [6, 8]$ we get

$$X^T X = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, \quad X^T \mathbf{y} = \begin{bmatrix} 14 \\ 6 \\ 8 \end{bmatrix}.$$

Solving the normal equations

The normal equations are

$$\begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \begin{bmatrix} b_0 \\ b_1 \\ b_2 \end{bmatrix} = \begin{bmatrix} 14 \\ 6 \\ 8 \end{bmatrix}.$$

Since the first column of $X^T X$ is the sum of the next two, $X^T X$ is not of full rank. It is easy to see that $r(X^T X) = 2$.

Solving the normal equations

To find a conditional inverse of $X^T X$, we apply Theorem 6.2,
using the nonsingular minor $\begin{bmatrix} 2 & 1 \\ 1 & 1 \end{bmatrix}$.

This gives

$$(X^T X)^c = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

and therefore

$$\mathbf{b} = (X^T X)^c X^T \mathbf{y} = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 14 \\ 6 \\ 8 \end{bmatrix} = \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix}.$$

one way classification
 simple structure
 J
 diagonal conditional inverse
 for $X^T X$

Solving the normal equations

However, using the minor $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ gives the conditional inverse

$$(X^T X)^c = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix},$$

which gives the solution

$$\mathbf{b} = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 14 \\ 6 \\ 8 \end{bmatrix} = \begin{bmatrix} 0 \\ 6 \\ 8 \end{bmatrix}.$$

Both these solutions solve the normal equations, and are equally valid! This is the problem with the less than full rank model.

Carbon removal example

Example. Consider the earlier carbon removal example. We have

$$X^T X = \begin{bmatrix} 9 & 3 & 3 & 3 \\ 3 & 3 & 0 & 0 \\ 3 & 0 & 3 & 0 \\ 3 & 0 & 0 & 3 \end{bmatrix}$$

so a conditional inverse is

$$(X^T X)^c = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{3} & 0 & 0 \\ 0 & 0 & \frac{1}{3} & 0 \\ 0 & 0 & 0 & \frac{1}{3} \end{bmatrix}.$$

Carbon removal example

We can also calculate

$$X^T \mathbf{y} = \begin{bmatrix} 303.3 \\ 105 \\ 117.9 \\ 80.4 \end{bmatrix}.$$

Using the conditional inverse above gives us a solution to the normal equations:

$$\mathbf{b} = (X^T X)^c X^T \mathbf{y} = \begin{bmatrix} 0 \\ 35 \\ 39.3 \\ 26.8 \end{bmatrix}.$$

Solving the normal equations

If the model is less than full rank, the normal equations have an infinite number of solutions.

Theorem 6.6

Let $A\mathbf{x} = \mathbf{g}$ be a consistent system. Then

$$\mathbf{x} = A^c \mathbf{g} + (I - A^c A) \mathbf{z}$$

solves the system, where \mathbf{z} is an arbitrary $p \times 1$ vector.

Solving the normal equations

Proof. We know that $A^c\mathbf{g}$ solves the system, so

$$\begin{aligned} A\mathbf{x} &= A[A^c\mathbf{g} + (I - A^cA)\mathbf{z}] \\ &= AA^c\mathbf{g} + (A - AA^cA)\mathbf{z} \\ &= \mathbf{g} + (A - A)\mathbf{z} = \mathbf{g}. \end{aligned}$$

Thus, for the normal equations, any vector of the form

$$\mathbf{b} = (X^T X)^c X^T \mathbf{y} + [I - (X^T X)^c X^T X]\mathbf{z}$$

satisfies the equations.

Solving the normal equations

Example. In the two-class example above, one solution to the normal equations was $(X^T X)^c X^T \mathbf{y} = [\begin{array}{ccc} 8 & -2 & 0 \end{array}]^T$.

Using the same conditional inverse we have

$$(X^T X)^c X^T X = \left[\begin{array}{ccc} 1 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{array} \right] \left[\begin{array}{ccc} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{array} \right] = \left[\begin{array}{ccc} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{array} \right].$$

Solving the normal equations

Then another solution to the normal equations is

$$\begin{aligned}
 \mathbf{b} &= (X^T X)^c X^T \mathbf{y} + [I - (X^T X)^c X^T X] \mathbf{z} \\
 &= \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix} + \left(\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{bmatrix} \right) \begin{bmatrix} z_1 \\ z_2 \\ z_3 \end{bmatrix} \\
 &= \begin{bmatrix} 8 - z_3 \\ -2 + z_3 \\ z_3 \end{bmatrix}
 \end{aligned}$$

for arbitrary z_3 .

For example, $\begin{bmatrix} 7 \\ -1 \\ 1 \end{bmatrix}$ is a solution.

Solving the normal equations

The converse of the above theorem is also true: *all* solutions to the system can be expressed in this form.

Theorem 6.7

Let $Ax = g$ be a consistent system and let x_0 be any solution to the system. Then for any A^c ,

$$\mathbf{x}_0 = A^c \mathbf{g} + (I - A^c A) \mathbf{z}$$

where $\mathbf{z} = \mathbf{x}_0$.

Solving the normal equations

Proof. Since \mathbf{x}_0 solves the system, we have

$$\begin{aligned} A^c \mathbf{g} + (I - A^c A) \mathbf{z} &= A^c \mathbf{g} + (I - A^c A) \mathbf{x}_0 \\ &= A^c \mathbf{g} + \mathbf{x}_0 - A^c A \mathbf{x}_0 \\ &= A^c \mathbf{g} + \mathbf{x}_0 - A^c \mathbf{g} = \mathbf{x}_0. \end{aligned}$$

For the normal equations, this means that any solution can be expressed as

$$\mathbf{b} = (X^T X)^c X^T \mathbf{y} + [I - (X^T X)^c X^T X] \mathbf{z}$$

for *any* conditional inverse $(X^T X)^c$, and some \mathbf{z} .

Solving the normal equations

Example. In the two-class example, we found the solution

$$\mathbf{b}_1 = \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix}$$

using our original conditional inverse.

But we also noted that the conditional inverse

$$(X^T X)_2^c = \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix}$$

produces the solution

$$\mathbf{b}_2 = \begin{bmatrix} 0 \\ 6 \\ 8 \end{bmatrix}.$$

Solving the normal equations

Using the theorem, the first solution can be written in terms of the second solution:

$$\begin{aligned}
 \mathbf{b}_1 &= (X^T X)_2^c X^T \mathbf{y} + (I - (X^T X)_2^c X^T X) \mathbf{z} \\
 &= \begin{bmatrix} 0 \\ 6 \\ 8 \end{bmatrix} + \left(\begin{bmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} - \begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \right) \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} 0 \\ 6 \\ 8 \end{bmatrix} + \begin{bmatrix} 1 & 0 & 0 \\ -1 & 0 & 0 \\ -1 & 0 & 0 \end{bmatrix} \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix} \\
 &= \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix}.
 \end{aligned}$$

Exam marks example

We compare the marks of students in 3 different mathematics classes. There is another factor (IQ), but we ignore this for the time being.

```
> maths <- read.csv("../data/math.csv")
> str(maths)

'data.frame':      30 obs. of  5 variables:
 $ X      : int  1 2 3 4 5 6 7 8 9 10 ...
 $ maths.y: int  81 84 81 79 78 79 81 85 72 79 ...
 $ iq     : int  99 103 108 109 96 104 96 105 94 91 ...
 $ class  : int  1 1 1 1 1 1 1 1 1 1 ...
 $ class.f: int  1 1 1 1 1 1 1 1 1 1 ...
```

*numerically : class 2 to class 1
is different from class 3 to
class 2*

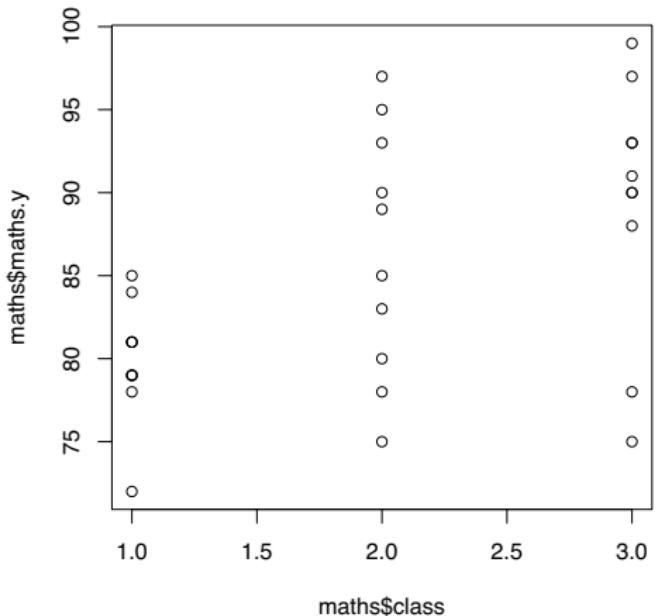
```
> maths$class.f <- factor(maths$class.f)
```

categorical but not ordinal

*categorically: class 2 to class 1
is the same as class 3 to 2.*

Exam marks example

```
> plot(maths$class, maths$maths.y)
```



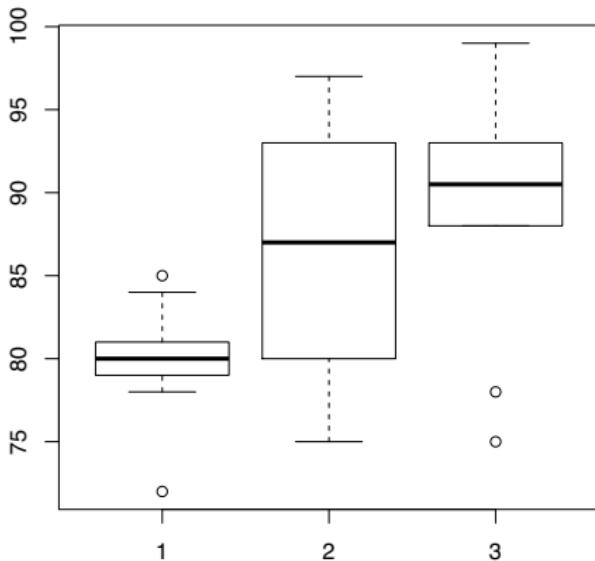
scatter plot

what R would do with 2 [numerical value]

!

Exam marks example

```
> plot(maths$class.f, maths$maths.y)
```



Exam marks example

```
> (y <- maths$maths.y)
[1] 81 84 81 79 78 79 81 85 72 79 85 78 93 80 83 95 90 89 97 75 90 75 99 97 93
[26] 91 88 93 90 78

> n <- dim(maths)[1]
> k <- length(levels(maths$class.f))
> X <- matrix(0,n,k+1)
> X[,1] <- 1
> X[maths$class.f==1,2] <- 1
 $\nearrow \text{In second column, all rows where class=1 set to be 1}$ 
> X[maths$class.f==2,3] <- 1
> X[maths$class.f==3,4] <- 1
```

Exam marks example

```
> X      mu class1 class2 class3
      [,1] [,2] [,3] [,4]
[1,]    1    1    0    0
[2,]    1    1    0    0
[3,]    1    1    0    0
[4,]    1    1    0    0
[5,]    1    1    0    0
[6,]    1    1    0    0
[7,]    1    1    0    0
[8,]    1    1    0    0
[9,]    1    1    0    0
[10,]   1    1    0    0
[11,]   1    0    1    0
[12,]   1    0    1    0
[13,]   1    0    1    0
[14,]   1    0    1    0
[15,]   1    0    1    0
[16,]   1    0    1    0
[17,]   1    0    1    0
[18,]   1    0    1    0
[19,]   1    0    1    0
[20,]   1    0    1    0
```

Exam marks example: reparametrisation

$$y_{ij} = \mu + \tau_i + \varepsilon_{ij}$$

```
> Xre <- X[,-1]
> (b <- solve(t(Xre) %*% Xre, t(Xre) %*% y))
      [,1]
[1,] 79.9
[2,] 86.5
[3,] 89.4
```

↳ reparametrise
get full rank model.

Exam marks example: reparametrisation

```
> modelre <- lm(y ~ 0 + X[,2] + X[,3] + X[,4])
> summary(modelre)
```

Call:

```
lm(formula = y ~ 0 + X[, 2] + X[, 3] + X[, 4])
```

Residuals:

Min	1Q	Median	3Q	Max
-14.40	-1.80	0.85	3.60	10.50

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
X[, 2]	79.900	2.053	38.92	<2e-16 ***
X[, 3]	86.500	2.053	42.14	<2e-16 ***
X[, 4]	89.400	2.053	43.55	<2e-16 ***

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 6.492 on 27 degrees of freedom

Multiple R-squared: 0.9948. Adjusted R-squared: 0.9942

Exam marks example

Let's look at the normal equations.

```
> t(X) %*% X
```

$$x^T x$$

```
 [,1] [,2] [,3] [,4]  
[1,] 30 10 10 10  
[2,] 10 10 0 0  
[3,] 10 0 10 0  
[4,] 10 0 0 10
```

```
> t(X) %*% y
```

$$x^T y$$

```
 [,1]  
[1,] 2558  
[2,] 799  
[3,] 865  
[4,] 894
```

Exam marks example

```
> XtXc <- matrix(0,4,4)
> XtXc[2:4,2:4] <- solve((t(X) %*% X)[2:4,2:4])
> (b <- XtXc %*% t(X) %*% y)

[,1]
[1,] 0.0
[2,] 79.9
[3,] 86.5
[4,] 89.4

> round(t(X) %*% X %*% b - t(X) %*% y, 3)

[,1]
[1,] 0
[2,] 0
[3,] 0
[4,] 0
```

Exam marks example

```
> (b2 <- ginv(t(X) %*% X) %*% t(X) %*% y)
[1,] [,1]
[1,] 63.95
[2,] 15.95
[3,] 22.55
[4,] 25.45

> round(t(X) %*% X %*% b2 - t(X) %*% y, 3)
[1,] [,1]
[1,] 0
[2,] 0
[3,] 0
[4,] 0
```

Exam marks example

```
> I4 <- diag(4)
> z <- c(2,8,-2,1)
> (b3 <- b + (I4 - XtXc %*% t(X) %*% X) %*% z)

[,1]
[1,] 2.0
[2,] 77.9
[3,] 84.5
[4,] 87.4

> round(t(X) %*% X %*% b3 - t(X) %*% y, 3)

[,1]
[1,] 0
[2,] 0
[3,] 0
[4,] 0
```

Exam marks example

```
> b + (I4 - XtXc %*% t(X) %*% X) %*% b3
```

```
[,1]
[1,] 2.0
[2,] 77.9
[3,] 84.5
[4,] 87.4
```

$$b + (I - (X^T X)^{-1} X^T X) b^3$$

```
> b3
```

```
[,1]
[1,] 2.0
[2,] 77.9
[3,] 84.5
[4,] 87.4
```

Estimability

Now we know how to solve the normal equations; furthermore, we know how to find *all* solutions for them.

But which solution(s) do we want?

Or rather, which solutions can we find?

Estimability

Some quantities do not change no matter what solutions we use for the normal equations. We call these quantities estimable.

A trivial example is the responses \mathbf{y} .

Definition 6.8

In the general linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, a function $\mathbf{t}^T\boldsymbol{\beta}$ is said to be estimable if there exists a vector \mathbf{c} such that $E[\mathbf{c}^T\mathbf{y}] = \mathbf{t}^T\boldsymbol{\beta}$.

In other words, a quantity is estimable if there is a linear unbiased estimator for it.

$\mathbf{t}^T\boldsymbol{\beta}$ is a linear combination of \mathbf{y} .
 $\mathbf{t}^T\boldsymbol{\beta}$ is a linear unbiased estimator of $\mathbf{t}^T\boldsymbol{\beta}$.

exp of (comb of response).
is comb of parameter
responses are known,
didn't depend on
parametrisation

Estimability

*Estimator of $t^T \beta$
 $\hat{\theta} = t^T b$*

Theorem 6.9

In the general linear model $\mathbf{y} = X\beta + \epsilon$, $t^T \beta$ is estimable if and only if there is a solution to the linear system $X^T X z = t$.

Proof. (\Leftarrow) Let \mathbf{z}_0 be a solution to $X^T X \mathbf{z} = \mathbf{t}$ and put $\mathbf{c} = X \mathbf{z}_0$.

Then

$$E[\mathbf{c}^T \mathbf{y}] = E[\mathbf{z}_0^T X^T \mathbf{y}] = \mathbf{z}_0^T X^T E[\mathbf{y}] = \mathbf{z}_0^T X^T X \beta = \mathbf{t}^T \beta,$$

so $t^T \beta$ is estimable.

Estimability

Example. Consider our two-class example. We had

$$X = \begin{bmatrix} 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}, \quad X^T X = \begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix}.$$

Consider the combination of parameters $\beta_1 - \beta_2$. This corresponds to $t^T \beta$ where

$$t = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}.$$

Estimability

We look for a solution to the system

$$\begin{matrix} & \begin{matrix} z_1 \\ z_2 \\ z_3 \end{matrix} \\ \begin{matrix} x_1 \\ x_2 \\ x_3 \end{matrix} & \begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \end{matrix} = \begin{bmatrix} 0 \\ 1 \\ -1 \end{bmatrix}.$$

A

$\left[\begin{matrix} 1 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{matrix} \right] - \left[\begin{matrix} 0 & 0 & 0 \\ 0 & 0 & 0 \\ 0 & 0 & 0 \end{matrix} \right]$

$\text{rank}(A(g)) = \text{rank}(A)$.
consistent.

This system has solution $z_1 = 0, z_2 = 1, z_3 = -1$, so $\beta_1 - \beta_2$ is estimable.

Estimability

estimability is a property of model
(no response inside)

Theorem 6.10

In the general linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, $\mathbf{t}^T\boldsymbol{\beta}$ is estimable if and only if

$$\mathbf{t}^T(\mathbf{X}^T\mathbf{X})^c\mathbf{X}^T\mathbf{X} = \mathbf{t}^T,$$

for some (and thus all) conditional inverse of $(\mathbf{X}^T\mathbf{X})$.

Proof. (\Leftarrow) Assume that $\mathbf{t}^T(\mathbf{X}^T\mathbf{X})^c\mathbf{X}^T\mathbf{X} = \mathbf{t}^T$, so

take transpose of both sides

$$\mathbf{X}^T\mathbf{X}((\mathbf{X}^T\mathbf{X})^c)^T\mathbf{t} = \mathbf{X}^T\mathbf{X}(\mathbf{X}^T\mathbf{X})^c\mathbf{t} = \mathbf{t}.$$

$$\mathbf{X}^T\mathbf{X}\mathbf{z} = \mathbf{t}$$

$$\text{let } \mathbf{z} = (\mathbf{X}^T\mathbf{X})^c\mathbf{t}$$

This means that $(\mathbf{X}^T\mathbf{X})^c\mathbf{t}$ is a solution to the system $\mathbf{X}^T\mathbf{X}\mathbf{z} = \mathbf{t}$, and Theorem 6.9 implies that $\mathbf{t}^T\boldsymbol{\beta}$ is estimable.

Estimability

(\Rightarrow) Suppose that $\mathbf{t}^T \boldsymbol{\beta}$ is estimable. By Theorem 6.9, there exists a solution to the system $X^T X \mathbf{z} = \mathbf{t}$.

A solution for $Ax=g$ is $A^c g$

We know that a solution is $\mathbf{z} = (X^T X)^c \mathbf{t}$. (Note that the conditional inverse is arbitrary.)

In other words,

$$X^T X (X^T X)^c \mathbf{t} = \mathbf{t}$$

and by taking transposes, we see that this gives the required condition.

Estimability

Example. Consider the previous example. Let us take the conditional inverse

$$(X^T X)^c = \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix}$$

and consider again the quantity $\beta_1 - \beta_2$, which corresponds to
 $\mathbf{t} = [0 \ 1 \ -1]^T$.

Estimability

Then

$$\begin{aligned}
 \mathbf{t}^T (X^T X)^c (X^T X) &= [0 \ 1 \ -1] \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \\
 &= [0 \ 1 \ -1] \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{bmatrix} \\
 &= [0 \ 1 \ -1] = \mathbf{t}^T,
 \end{aligned}$$

so again we see that $\beta_1 - \beta_2$ is estimable.

Estimability

On the other hand, suppose we take $\mathbf{t} = [1 \ 0 \ 0]^T$ so that $\mathbf{t}^T \boldsymbol{\beta} = \beta_0$.

Then we have

$$\begin{aligned}\mathbf{t}^T (X^T X)^c (X^T X) &= [1 \ 0 \ 0] \begin{bmatrix} 1 & -1 & 0 \\ -1 & 2 & 0 \\ 0 & 0 & 0 \end{bmatrix} \begin{bmatrix} 2 & 1 & 1 \\ 1 & 1 & 0 \\ 1 & 0 & 1 \end{bmatrix} \\ &= [1 \ 0 \ 0] \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & -1 \\ 0 & 0 & 0 \end{bmatrix} \\ &= [1 \ 0 \ 1] \neq \mathbf{t}^T,\end{aligned}$$

so β_0 is not estimable.

Estimability

Example. We return to the carbon removal example. We are interested in seeing if the three carbon removal treatments have (significantly) different means.

To test this, we look at the quantities $\tau_1 - \tau_2$ and $\tau_1 - \tau_3$.

If both of these are (close to) 0, then the treatments are not significantly different.

Estimability

We have

$$X^T X = 3 \begin{bmatrix} 3 & 1 & 1 & 1 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}, \quad (X^T X)^c X^T X = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix}$$

minor

and the coefficient vectors

$$\mathbf{t}_1 = \begin{bmatrix} \mu \\ t_1 \\ t_2 \\ t_3 \end{bmatrix} = \begin{bmatrix} 0 \\ 1 \\ -1 \\ 0 \end{bmatrix}, \quad \mathbf{t}_2 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ -1 \end{bmatrix}.$$

Estimability

$$\mathbf{t}_1^T (X^T X)^c X^T X = \begin{bmatrix} 0 & 1 & -1 & 0 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & -1 & 0 \end{bmatrix}$$

so $\mathbf{t}_1^T \boldsymbol{\beta} = \tau_1 - \tau_2$ is estimable.

$$\mathbf{t}_2^T (X^T X)^c X^T X = \begin{bmatrix} 0 & 1 & 0 & -1 \end{bmatrix} \begin{bmatrix} 0 & 0 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & -1 \end{bmatrix}$$

so $\mathbf{t}_2^T \boldsymbol{\beta} = \tau_1 - \tau_3$ is also estimable.

Estimability

full rank

$$\underline{t^T \beta} \Rightarrow \frac{\underline{t^T b}}{\underline{t^T (X^T X)^{-1} X^T y}}$$

Next we will prove that no matter what conditional inverse we use, we will still generate the same estimate for an estimable quantity.

less than full rank

$$\begin{aligned} \underline{t^T \beta} &\Rightarrow \underline{z^T X^T y} \\ &= \underline{[(X^T X)^C t]^T X^T y} \\ &= \underline{t^T (X^T X)^C X^T y} \\ &= \underline{t^T b} \end{aligned}$$

Theorem 6.11 (A Gauss-Markov Theorem)

In the general linear model $y = X\beta + \epsilon$, suppose $t^T \beta$ is estimable. Then the best linear unbiased estimator (BLUE) for $t^T \beta$ is $\underline{z^T X^T y}$, where \underline{z} is a solution to the system $X^T X z = t$. Furthermore, this estimate is the same for any solution of the system, and can be written $\underline{t^T b}$, where \underline{b} is any solution to the normal equations.

BLUE should be in the format \underline{Ly}
 \downarrow
 unique

$$\begin{aligned} z &= (X^T X)^C t \\ z^T &= t^T (X^T X)^C X^T y \end{aligned}$$

$$\begin{aligned} X^T X b &= X^T y \\ b &= (X^T X)^C X^T y \end{aligned}$$

Estimability

Proof. We first show unbiasedness of the estimator.

$$\begin{aligned} E[\mathbf{z}^T X^T \mathbf{y}] &= \mathbf{z}^T X^T E[\mathbf{y}] \\ &= \mathbf{z}^T X^T X \boldsymbol{\beta} \\ &= \mathbf{t}^T \boldsymbol{\beta}. \end{aligned}$$

$$\mathbf{y} \sim N(\mathbf{x}\boldsymbol{\beta}, \sigma^2 \mathbf{I})$$

BLUEness is more involved, but is similar to the proof of Theorem 4.4.

Now suppose we have two solutions to the system $X^T X \mathbf{z} = \mathbf{t}$, called \mathbf{z}_0 and \mathbf{z}_1 . Let \mathbf{b} be any solution to the normal equations:

$$X^T X \mathbf{b} = X^T \mathbf{y}.$$

Estimability

The best linear unbiased estimator of $\mathbf{t}^T \boldsymbol{\beta}$ is

$$\mathbf{z}_0^T X^T \mathbf{y} = \mathbf{z}_0^T X^T X \mathbf{b} = (X^T X \mathbf{z}_0)^T \mathbf{b} = \mathbf{t}^T \mathbf{b}.$$

Since \mathbf{b} is an arbitrary solution, this is the same no matter what solution we choose.

Similarly,

$$\begin{aligned} & \text{so} \\ & \text{suppose } \mathbf{b}_1, \mathbf{b}_2 \\ & \left\{ \begin{array}{l} \mathbf{t}^T \mathbf{b}_1 = \mathbf{z}_0^T X^T \mathbf{y} \\ \mathbf{t}^T \mathbf{b}_2 = \mathbf{z}_0^T X^T \mathbf{y} \end{array} \right. \end{aligned}$$

$$\mathbf{z}_1^T X^T \mathbf{y} = \mathbf{t}^T \mathbf{b} = \mathbf{z}_0^T X^T \mathbf{y}.$$

Thus the best linear unbiased estimator is unique, and equal to $\underline{\mathbf{t}^T \mathbf{b}}$.

Estimability

Example. Let's look again at the two-class example. We know that $\beta_1 - \beta_2$ is estimable. We also know that solutions to the normal equations include

$$\mathbf{b} = \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix}, \quad \mathbf{b}_0 = \begin{bmatrix} 0 \\ 6 \\ 8 \end{bmatrix}.$$

To estimate $\beta_1 - \beta_2$, we can use

$$\mathbf{t}^T \mathbf{b} = [0 \ 1 \ -1] \begin{bmatrix} 8 \\ -2 \\ 0 \end{bmatrix} = -2.$$

Estimability

However, from Theorem 6.11, we can also use

$$\mathbf{t}^T \mathbf{b}_0 = [\begin{array}{ccc} 0 & 1 & -1 \end{array}] \left[\begin{array}{c} 0 \\ 6 \\ 8 \end{array} \right] = -2.$$

This estimate is the same as the previous one, which follows from the theorem: *any solution to the normal equation, using any conditional inverse, will produce exactly the same estimate.*

In other words, the estimator is unique.

Estimability

Example. Back to the carbon removal example. We have shown that $\tau_1 - \tau_2$ and $\tau_1 - \tau_3$ are estimable. We estimate them by

$$\mathbf{t}_1^T \mathbf{b} = [0 \ 1 \ -1 \ 0] \begin{bmatrix} 0 \\ 35 \\ 39.3 \\ 26.8 \end{bmatrix} = -4.3$$

and

$$\mathbf{t}_2^T \mathbf{b} = [0 \ 1 \ 0 \ -1] \begin{bmatrix} 0 \\ 35 \\ 39.3 \\ 26.8 \end{bmatrix} = 8.2$$

respectively.

Again, no matter what conditional inverse we use, these estimates remain the same.

Estimability theorems

Theorem 6.12

In the linear model $\mathbf{y} = X\beta + \varepsilon$, elements of $X\beta$ are estimable.

Proof. We know that $E[\mathbf{y}] = X\beta$. Now take \mathbf{e}_i to be the i th standard basis vector.

We have

$$\begin{aligned}(X\beta)_i &= \mathbf{e}_i^T X\beta \\ &= \mathbf{e}_i^T E[\mathbf{y}] \\ &= E[\mathbf{e}_i^T \mathbf{y}]\end{aligned}$$

and so the i th element of $X\beta$ is estimable.

 $\left[\begin{array}{c} \vdots \\ \vdots \\ \vdots \\ \vdots \\ \vdots \end{array} \right] \quad \text{i-th position}$

\Rightarrow consider estimability of a scalar.
 \Rightarrow each element of $X\beta$

Estimability theorems

Example. Consider the carbon removal example. We have

$$X = \begin{bmatrix} 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \\ 1 & 0 & 0 & 1 \end{bmatrix}, \quad \beta = \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \\ \tau_3 \end{bmatrix}.$$

We know that we cannot estimate the parameter vector β , because it is not uniquely determined.

Estimability theorems

However, the real quantities of interest are the mean responses from the three treatments. These are:

$$\begin{aligned}\mu + \tau_1 &= [1 \ 1 \ 0 \ 0] \beta && \text{1st of } X\beta \\ \mu + \tau_2 &= [1 \ 0 \ 1 \ 0] \beta && \text{2nd of } X\beta \\ \mu + \tau_3 &= [1 \ 0 \ 0 \ 1] \beta && :\end{aligned}$$

and each of these are elements of $X\beta$. Therefore, they are estimable.

In a one-way classification model with any number of levels, $\mu + \tau_i$ is always estimable.

subpopulation mean
is always estimable
 \Rightarrow given a sample from that subpopulation

Estimability theorems

Theorem 6.13

Let $\mathbf{t}_1^T \boldsymbol{\beta}, \mathbf{t}_2^T \boldsymbol{\beta}, \dots, \mathbf{t}_k^T \boldsymbol{\beta}$ be estimable functions, and let

$$z = a_1 \mathbf{t}_1^T \boldsymbol{\beta} + a_2 \mathbf{t}_2^T \boldsymbol{\beta} + \dots + a_k \mathbf{t}_k^T \boldsymbol{\beta}.$$

linear comb of estimable functions is estimable

Then z is estimable, and the best linear unbiased estimator for z is

$$a_1 \mathbf{t}_1^T \mathbf{b} + a_2 \mathbf{t}_2^T \mathbf{b} + \dots + a_k \mathbf{t}_k^T \mathbf{b}.$$

Estimability theorems

Proof. By definition,

$$\underline{\mathbf{t}}^T \boldsymbol{\beta}$$

$$z = (a_1 \mathbf{t}_1 + a_2 \mathbf{t}_2 + \dots + a_k \mathbf{t}_k)^T \boldsymbol{\beta}.$$

From Theorem 6.10,

$$\begin{aligned}
 & (a_1 \mathbf{t}_1 + a_2 \mathbf{t}_2 + \dots + a_k \mathbf{t}_k)^T (X^T X)^c X^T X \\
 &= a_1 \mathbf{t}_1^T (X^T X)^c X^T X + a_2 \mathbf{t}_2^T (X^T X)^c X^T X + \dots + a_k \mathbf{t}_k^T (X^T X)^c X^T X \\
 &= a_1 \mathbf{t}_1^T + a_2 \mathbf{t}_2^T + \dots + a_k \mathbf{t}_k^T \\
 &= (a_1 \mathbf{t}_1 + a_2 \mathbf{t}_2 + \dots + a_k \mathbf{t}_k)^T.
 \end{aligned}$$

Therefore z is estimable, with BLUE

$$(a_1 \mathbf{t}_1 + a_2 \mathbf{t}_2 + \dots + a_k \mathbf{t}_k)^T \mathbf{b}.$$

Estimability theorems

Of particular interest in many studies is the way different populations compare against each other. To attach a numerical value to these comparisons, we form linear combinations

$$a_1\tau_1 + a_2\tau_2 + \dots + a_k\tau_k,$$

where $\sum_{i=1}^k a_i = 0$.

These *treatment contrasts* wipe out the effect of the overall mean response, to describe the differences between populations.

Estimability theorems

In a one-way classification model, any treatment contrast is estimable.

If

$$z = a_1\tau_1 + a_2\tau_2 + \dots + a_k\tau_k$$

is a treatment contrast, then

$$\begin{aligned} z &= \sum_{i=1}^k a_i \mu + a_1\tau_1 + a_2\tau_2 + \dots + a_k\tau_k \\ &= a_1(\mu + \tau_1) + a_2(\mu + \tau_2) + \dots + a_k(\mu + \tau_k) \end{aligned}$$

*element of \mathcal{X}^**
 is a linear combination of the estimable functions $\mu + \tau_i$, and is therefore estimable.

Estimability theorems

Of particular interest among treatment contrasts is the contrast of the form $\tau_i - \tau_j$, for some $i \neq j$. This is because

$$\tau_i - \tau_j = (\mu + \tau_i) - (\mu + \tau_j)$$

is the difference between the mean responses in populations i and j .

We would expect to estimate this contrast by the corresponding difference in sample means, $\bar{y}_i - \bar{y}_j$. We can show using the theory we have developed that this is in fact the case.

Estimability theorems

Example. We do this for $k = 3$ and the contrast $\tau_1 - \tau_2$. Our matrices are

$$\mathbf{y} = \begin{bmatrix} y_{11} \\ \vdots \\ y_{1n_1} \\ y_{21} \\ \vdots \\ y_{2n_2} \\ y_{31} \\ \vdots \\ y_{3n_3} \end{bmatrix}, \quad X = \begin{bmatrix} 1 & 1 & 0 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 1 & 0 \\ 1 & 0 & 0 & 1 \\ \vdots & \vdots & \vdots & \vdots \\ 1 & 0 & 0 & 1 \end{bmatrix}, \quad \boldsymbol{\beta} = \begin{bmatrix} \mu \\ \tau_1 \\ \tau_2 \\ \tau_3 \end{bmatrix}.$$

Estimability theorems

Direct multiplication gives

$$X^T \mathbf{y} = \begin{bmatrix} \sum_{i=1}^3 \sum_{j=1}^{n_i} y_{ij} \\ \sum_j y_{1j} \\ \sum_j y_{2j} \\ \sum_j y_{3j} \end{bmatrix}, \quad X^T X = \begin{bmatrix} n & n_1 & n_2 & n_3 \\ n_1 & n_1 & 0 & 0 \\ n_2 & 0 & n_2 & 0 \\ n_3 & 0 & 0 & n_3 \end{bmatrix}.$$

total number
 ↓ # 1 ↓ # 2 ↓ # 3.

Using the conditional inverse algorithm on the lower right corner of $X^T X$ gives

$$(X^T X)^c = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{n_1} & 0 & 0 \\ 0 & 0 & \frac{1}{n_2} & 0 \\ 0 & 0 & 0 & \frac{1}{n_3} \end{bmatrix}.$$

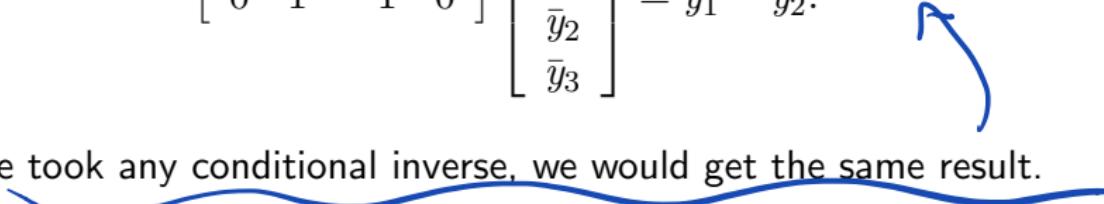
Estimability theorems

Therefore a solution to the normal equations is

$$\mathbf{b} = (X^T X)^c X^T \mathbf{y} = \begin{bmatrix} 0 \\ \bar{y}_1 \\ \bar{y}_2 \\ \bar{y}_3 \end{bmatrix}.$$

*nice formula for b
since we take lower right corner of $X^T X$*

We have $\tau_1 - \tau_2 = [0 \ 1 \ -1 \ 0] \boldsymbol{\beta}$, so the best linear unbiased estimator for $\tau_1 - \tau_2$ is

$$[0 \ 1 \ -1 \ 0] \begin{bmatrix} 0 \\ \bar{y}_1 \\ \bar{y}_2 \\ \bar{y}_3 \end{bmatrix} = \bar{y}_1 - \bar{y}_2.$$


If we took any conditional inverse, we would get the same result.

Exam marks example

We return to the `maths` dataset. Recall that b , b_2 and b_3 are all solutions to the normal equations.

```
> (tt <- c(0,1,-1,0))
[1] 0 1 -1 0
> round(tt %*% Xtxc %*% t(X) %*% X, 5) # estimable
[,1] [,2] [,3] [,4]
[1,] 0 1 -1 0
> (tt2 <- c(1,1,1,1)) μ + τ₁ + τ₂ + τ₃
[1] 1 1 1 1
> tt2 %*% Xtxc %*% t(X) %*% X # not estimable
[,1] [,2] [,3] [,4]
[1,] 3 1 1 1
```

Exam marks example

```
> tt %*% b
[1,]
[1,] -6.6

> tt %*% b2
[1,]
[1,] -6.6

> tt %*% b3
[1,]
[1,] -6.6

> mean(maths$maths.y[maths$class.f==1]) -
+      mean(maths$maths.y[maths$class.f==2])
[1] -6.6
```

Exam marks example

```
> tt2 %*% b  
      [,1]  
[1,] 255.8  
  
> tt2 %*% b2  
      [,1]  
[1,] 127.9  
  
> tt2 %*% b3  
      [,1]  
[1,] 251.8
```

Exam marks example

For the less than full rank model, R uses contrasts for its tests.

The two main contrast sets are `contr.treatment` and `contr.sum`. For the one-way classification model:

Label	<code>contr.treatment</code>	<code>contr.sum</code>
Intercept	μ_1	$\bar{\mu}$
factor1		$\mu_1 - \bar{\mu}$
factor2	$\mu_2 - \mu_1$	$\mu_2 - \bar{\mu}$
factor3	$\mu_3 - \mu_1$	$\mu_3 - \bar{\mu}$
:	:	:
factor($k-1$)	$\mu_{k-1} - \mu_1$	$\mu_{k-1} - \bar{\mu}$
factor(k)	$\mu_k - \mu_1$	

$$\mu_1 + \dots + \mu_k = k\bar{\mu}$$

$$\mu_k = k\bar{\mu} - (\mu_1 + \dots + \mu_{k-1})$$

$$\mu_k - \bar{\mu} = \underline{(k-1)\bar{\mu}} - \underline{(\mu_1 + \dots + \mu_{k-1})}$$

has to make $\bar{\mu}$?

$\bar{\mu}$, so it is fully determined

Exam marks example

$$\text{our model } y_i = \mu + \tau_i + \varepsilon_i$$

$$\text{level 1: } \mu + \tau_i$$

In terms of our parameters:

Label	contr.treatment	contr.sum	
Intercept	$\mu + \tau_1$	$\mu + \frac{1}{k} \sum \tau_i$	$\rightarrow \frac{1}{k} (\mu + \tau_1 + \mu + \tau_2 + \dots + \mu + \tau_k)$
factor1		$\tau_1 - \frac{1}{k} \sum \tau_i$	$\rightarrow \mu + \tau_1 - \mu - \frac{1}{k} \sum \tau_i$
factor2	$\tau_2 - \tau_1$	$\tau_2 - \frac{1}{k} \sum \tau_i$	
factor3	$\tau_3 - \tau_1$	$\tau_3 - \frac{1}{k} \sum \tau_i$	
\vdots	\vdots	\vdots	
factor($k-1$)	$\tau_{k-1} - \tau_1$	$\tau_{k-1} - \frac{1}{k} \sum \tau_i$	
factor(k)	$\tau_k - \tau_1$		

\rightarrow not weighted average, just average

Exam marks example

```
> contrasts(maths$class.f) <- contr.treatment(k)
> model <- lm(maths.y ~ class.f, data = maths)
> summary(model)
Call:
lm(formula = maths.y ~ class.f, data = maths)
```

what type of contrast you want to do
default contrast treatment

\downarrow **factor!**

Residuals:

Min	1Q	Median	3Q	Max
-14.40	-1.80	0.85	3.60	10.50

$$\begin{aligned} \mu + \tau_2 &= 79.9 + 6.6 \\ \mu + \tau_3 &= 79.9 + 9.5 \end{aligned}$$

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)	H ₀
(Intercept)	79.900	2.053	38.922	< 2e-16 ***	" $\mu + \tau_1 = 0$? "
class.f2	6.600	2.903	2.273	0.03117 *	" $\tau_2 - \tau_1 = \tau_2$? "
class.f3	9.500	2.903	3.272	0.00292 **	" $\tau_3 - \tau_1 = \tau_3$? "

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

Residual standard error: 6.492 on 27 degrees of freedom

Exam marks example

```
> contrasts(maths$class.f) <- contr.sum(k)
> model2 <- lm(maths.y ~ class.f, data = maths)
> summary(model2)
```

Call:

```
lm(formula = maths.y ~ class.f, data = maths)
```

Residuals:

Min	1Q	Median	3Q	Max
-14.40	-1.80	0.85	3.60	10.50

Coefficients:

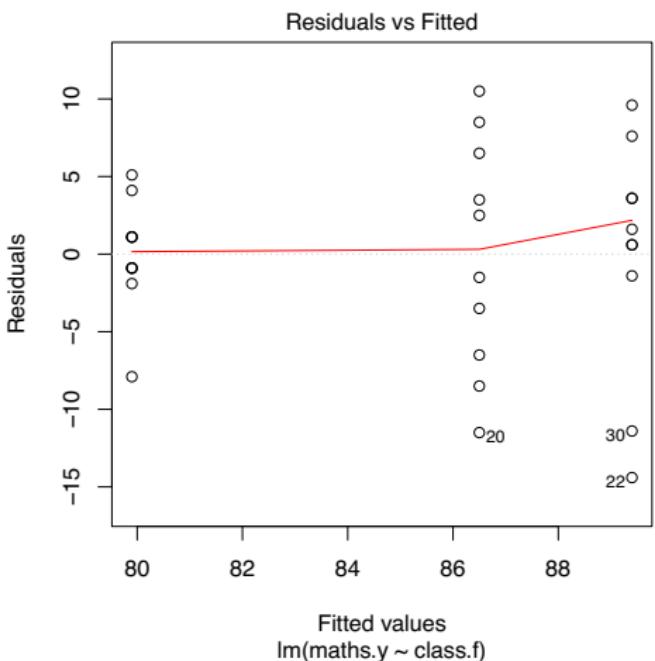
	Estimate	Std. Error	t value	Pr(> t)
(Intercept)	85.267	1.185	71.943	< 2e-16 ***
class.f1	-5.367	1.676	-3.202	0.00348 **
class.f2	1.233	1.676	0.736	0.46818

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

Residual standard error: 6.492 on 27 degrees of freedom

Exam marks example

```
> plot(model, which=1)
```



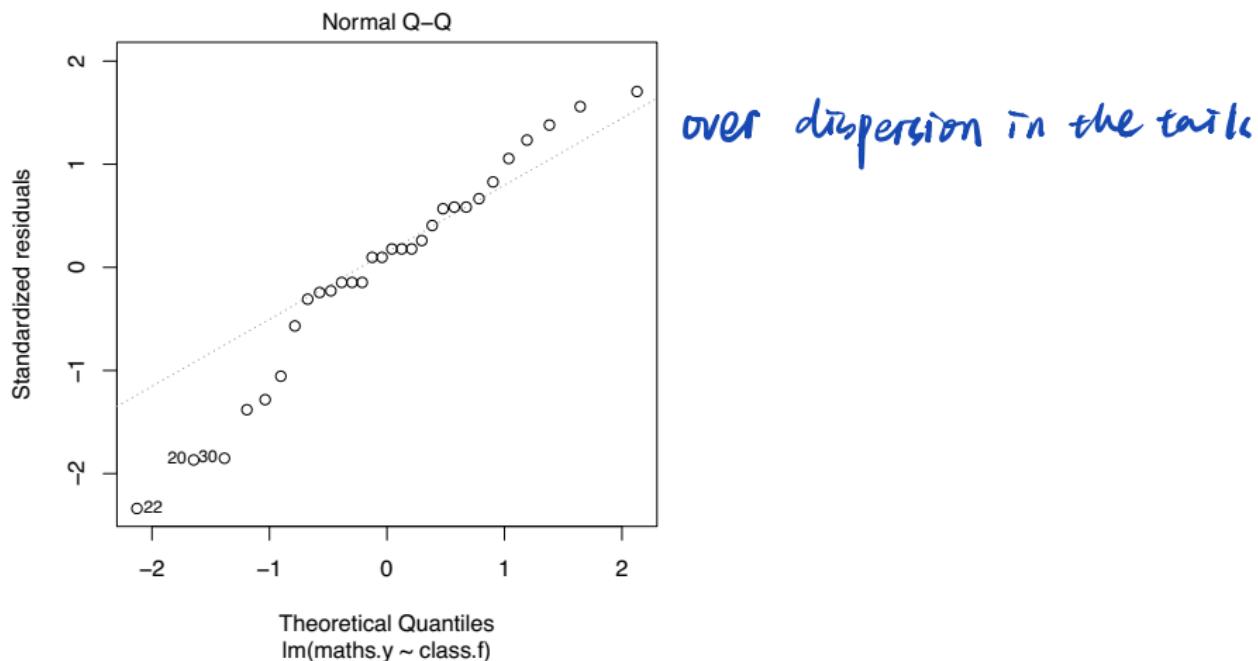
?

3 fitted value

since 3 different classes.

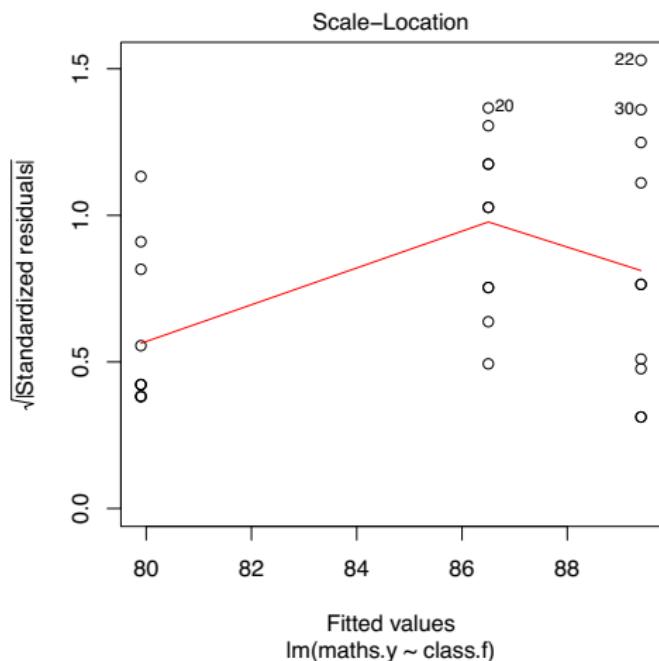
Exam marks example

```
> plot(model, which=2)
```



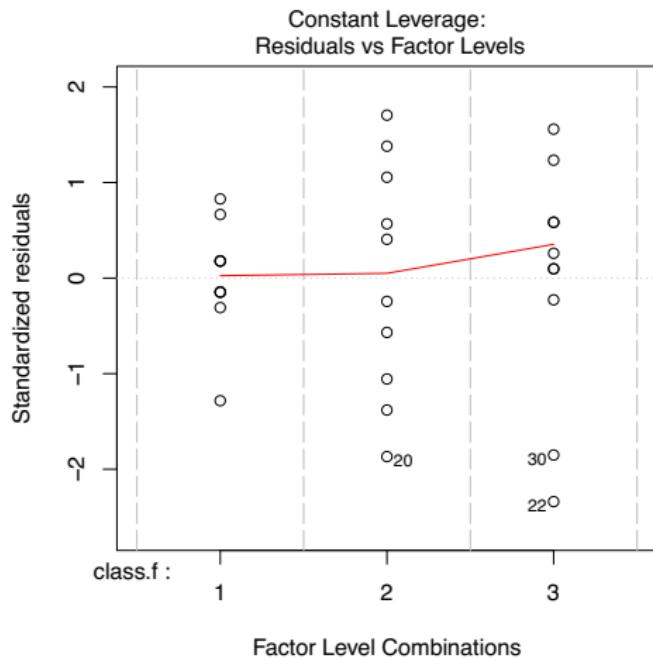
Exam marks example

```
> plot(model, which=3)
```



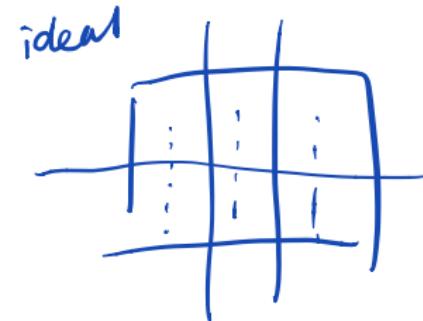
Exam marks example

> *plot(model, which=5)*



constant leverage. ?

one-way classification



Estimating σ^2 in the less than full rank model

In the full rank model, we estimated σ^2 by

$$s^2 = \frac{SS_{Res}}{n - p},$$

where n is the sample size, p is the number of parameters, and SS_{Res} is the sum of squares of the residuals:

$$SS_{Res} = (\mathbf{y} - X\mathbf{b})^T(\mathbf{y} - X\mathbf{b}) = \mathbf{y}^T[I - X(X^T X)^{-1}X^T]\mathbf{y}.$$

Estimating σ^2 in the less than full rank model

For the less than full rank model we can still define the residual sum of squares as

$$SS_{Res} = (\mathbf{y} - X\mathbf{b})^T(\mathbf{y} - X\mathbf{b}),$$

where \mathbf{b} is any solution to the normal equations.

Although \mathbf{b} can vary, $X\mathbf{b}$ will not, because $X\beta$ is estimable.
Therefore SS_{Res} is invariant to the choice of \mathbf{b} .

Estimating σ^2 in the less than full rank model

Theorem 6.14

$$SS_{Res} = \mathbf{y}^T [I - X(X^T X)^c X^T] \mathbf{y}.$$

Proof. Let $\mathbf{b} = (X^T X)^c X^T \mathbf{y}$. Then

$$\begin{aligned}
 SS_{Res} &= (\mathbf{y}^T - \mathbf{b}^T X^T)(\mathbf{y} - X\mathbf{b}) \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X\mathbf{b} + \mathbf{b}^T X^T X\mathbf{b} \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X(X^T X)^c X^T \mathbf{y} + \mathbf{y}^T X(X^T X)^c X^T X(X^T X)^c X^T \mathbf{y} \\
 &= \mathbf{y}^T \mathbf{y} - 2\mathbf{y}^T X(X^T X)^c X^T \mathbf{y} + \mathbf{y}^T X(X^T X)^c X^T \mathbf{y} \\
 &= \mathbf{y}^T [I - \underbrace{X(X^T X)^c X^T}_\text{Hat matrix}] \mathbf{y}.
 \end{aligned}$$

$$\mathbf{A}^c \mathbf{A} \mathbf{A}^c \neq \mathbf{A}^c$$

$$\mathbf{A} \mathbf{A}^c \mathbf{A} = \mathbf{A}$$

= \mathbf{X}

Estimating σ^2 in the less than full rank model

How do we find an estimator for σ^2 ?

$$\begin{aligned}
 SS_{Res} &= e^T e \\
 &= (y - X\beta)^T (y - X\beta) \\
 &= [(I - H)y]^T [(I - H)y] \\
 &= y^T (I - H)^T (I - H)y \\
 &\quad \text{idempotent} \\
 &\quad \text{symmetric} \\
 &\geq y^T (I - H)y
 \end{aligned}$$

Let's consider SS_{Res} again. Take $H = X(X^T X)^{-1}X^T$ and remember that $HX = X$.

$$y \sim N(X\beta, \sigma^2 I)$$

$$\begin{aligned}
 E[SS_{Res}] &= E[y^T (I - H)y] \\
 &= \text{tr}(I - H)\sigma^2 + (X\beta)^T (I - H)X\beta \\
 &= \text{tr}(I - H)\sigma^2 + \beta^T X^T X\beta - \beta^T X^T H X\beta \\
 &= \text{tr}(I - H)\sigma^2 + \cancel{\beta^T X^T X\beta} - \cancel{\beta^T X^T H X\beta} \\
 &= \text{tr}(I - H)\sigma^2.
 \end{aligned}$$

Estimating σ^2 in the less than full rank model

Since $I - H$ is symmetric and idempotent, we have

$$E[SS_{Res}] = r(I - H)\sigma^2 = (n - r)\sigma^2,$$

where $r = r(X)$, the rank of X .

Theorem 6.15

In the general linear model $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, suppose X has rank r and $\boldsymbol{\varepsilon}$ has mean $\mathbf{0}$ and variance $\sigma^2 I$. Then an unbiased estimator for σ^2 is

$$\frac{SS_{Res}}{n - r}.$$

not the # of parameters p.
since # of parameters is too many [over parametrizing]
true # of parameters p

Estimating σ^2 in the less than full rank model

Example. We return to the carbon removal example.

```
> y <- c(34.6,35.1,35.3,38.8,39.0,40.1,26.7,26.7,27.0)
> X <- matrix(c(rep(1,9),rep(0,27)),9,4)
> X[1:3,2] <- 1
> X[4:6,3] <- 1
> X[7:9,4] <- 1
> X
```

```
      [,1]  [,2]  [,3]  [,4]
[1,]    1    1    0    0
[2,]    1    1    0    0
[3,]    1    1    0    0
[4,]    1    0    1    0
[5,]    1    0    1    0
[6,]    1    0    1    0
[7,]    1    0    0    1
[8,]    .    .    .    .
[9,]    .    .    .    .
```

Estimating σ^2 in the less than full rank model

```
> (b <- ginv(t(X) %*% X) %*% t(X) %*% y)
[1,]
[1,] 25.275
[2,] 9.725
[3,] 14.025
[4,] 1.525

> e <- y - X %*% b
> (SSRes <- sum(e^2))

[1] 1.3

> (s2 <- SSRes/(9-3))

[1] 0.2166667
```

Exam marks example

```
> library(Matrix)
> (SSRes <- sum((y-X%*%b)^2))
[1] 1137.8
 $y^T X (X^T X)^{-1} X^T y$ 
> sum(y^2) - t(y) %*% X %*% XtXc %*% t(X) %*% y      ssres
[1,]
[1,] 1137.8
> (s2 <- SSRes/(n - rankMatrix(X)[1]))
[1] 42.14074
```

Exam marks example

```
> deviance(model)    sympy
[1] 1137.8
> deviance(model)/model$df.residual
[1] 42.14074
```

$$\frac{1}{n-r}$$

Interval estimation in the less than full rank model

We can find point estimates for estimable quantities. The next step is to try and find confidence intervals for them.

For the Gauss-Markov theorem we only required that ϵ has mean **0** and variance $\sigma^2 I$. However, to find confidence intervals, we need some idea of the distribution of the variables, so we suppose that $\epsilon \sim MVN(\mathbf{0}, \sigma^2 I)$.

Interval estimation in the less than full rank model

Recall that in the full rank model, we generated confidence intervals by finding a t -distributed quantity, which was created by dividing a normal variable by (the square root of) a χ^2 variable.

$$t \sim \mathcal{Z} / \sqrt{\chi^2 / r}$$

The χ^2 variable was

$$\frac{SS_{Res}}{\sigma^2},$$

which had $n - p$ degrees of freedom.

The σ^2 term was not known, but cancelled out another σ^2 term in the numerator to leave us with something that we could calculate.

We can proceed in a similar manner for the less than full rank model.

Interval estimation in the less than full rank model

Theorem 6.16

In the general linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, assume $\boldsymbol{\varepsilon} \sim MVN(\mathbf{0}, \sigma^2 I)$. Then

$$s^2 = \frac{SS_{Res}}{n-r}$$

$$\frac{(n-r)s^2}{\sigma^2} = \frac{SS_{Res}}{\sigma^2}$$

has a χ^2 distribution with $n - r$ degrees of freedom.

Theorem 6.17

In the general linear model $\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon}$, assume $\boldsymbol{\varepsilon} \sim MVN(\mathbf{0}, \sigma^2 I)$. If $\mathbf{t}^T \boldsymbol{\beta}$ is estimable, then $\mathbf{t}^T \mathbf{b}$ is independent of s^2 .

eg. in the full rank
 $\text{var}(\mathbf{b}) = (\mathbf{X}^T \mathbf{X})^{-1} \sigma^2$

& \mathbf{b} is independent of σ^2

Interval estimation in the less than full rank model

The steps to derive a confidence interval are very similar to that for the full rank case, but with two small differences. Firstly, we can only find confidence intervals for quantities that are estimable!

Secondly, we replace the inverse $(X^T X)^{-1}$ by the conditional inverse $(X^T X)^c$.

All other steps are the same.

Interval estimation in the less than full rank model

confidence interval for $t^T \beta$

We have

$$\begin{aligned} \text{Var } t^T b &= \text{Var } t^T (X^T X)^c X^T y \\ &= t^T (X^T X)^c X^T \sigma^2 I X (X^T X)^c t \\ &= \sigma^2 t^T (X^T X)^c t. \end{aligned}$$

$$\begin{aligned} E[t^T b] &= t^T \beta \\ \text{var}(Ay) &= A \text{var } y A^T \end{aligned}$$

If $t^T \beta$ is estimable

$$t^T (X^T X)^c (X^T X) = t^T.$$

Thus

$$\frac{(t^T b - t^T \beta) / \sigma \sqrt{t^T (X^T X)^c t}}{\sqrt{s^2 / \sigma^2}} \rightarrow \text{like } \frac{x - \bar{x}}{s}.$$

has a t distribution with $n - r$ degrees of freedom.

Interval estimation in the less than full rank model

This gives us the confidence interval for the (estimable) quantity $t^T \beta$, using a t distribution with $n - r$ degrees of freedom:

$$t^T b \pm t_{\alpha/2} s \sqrt{t^T (X^T X)^{-1} t}$$

↑
point
estimation

extra proof.
prove that $t^T (X^T X)^{-1} t$ is unique.

This formula can also be used to find confidence intervals for the individual parameters, if they are estimable.

$$t^T (X^T X)^{-1} (X^T X) = t^T$$

Interval estimation in the less than full rank model

Example. We return again to the carbon removal example.

Suppose we want to find a 95% confidence interval for $\tau_1 - \tau_2$.

```
> (tt <- c(0,1,-1,0))
[1] 0 1 -1 0
> ta <- qt(0.975,9-3) parameter  $t^T$   $(X^T X)^{-1}$   $t$ 
> halfwidth <- ta*sqrt(s2*t(tt)%*%ginv(t(X)%*%X)%*%tt)
> tt%*%b + c(-1,1)*halfwidth
[1] -5.22997 -3.37003
```

In particular, we can say with 95% confidence that the the first carbon removal treatment is not as effective as the second.

Interval estimation in the less than full rank model

Example. We showed earlier that in a 3-level 1-way classification model, the contrast $\tau_1 - \tau_2$ can be estimated by the difference in the respective ~~population~~
sample means, $\bar{y}_1 - \bar{y}_2$.

We also had

$$\mathbf{t} = \begin{bmatrix} 0 \\ 1 \\ -1 \\ 0 \end{bmatrix}, \quad (X^T X)^c = \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{n_1} & 0 & 0 \\ 0 & 0 & \frac{1}{n_2} & 0 \\ 0 & 0 & 0 & \frac{1}{n_3} \end{bmatrix}.$$

Interval estimation in the less than full rank model

Therefore we have

$$\begin{aligned} \mathbf{t}^T (X^T X)^c \mathbf{t} &= [0 \ 1 \ -1 \ 0] \begin{bmatrix} 0 & 0 & 0 & 0 \\ 0 & \frac{1}{n_1} & 0 & 0 \\ 0 & 0 & \frac{1}{n_2} & 0 \\ 0 & 0 & 0 & \frac{1}{n_3} \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ -1 \\ 0 \end{bmatrix} \\ &= \frac{1}{n_1} + \frac{1}{n_2} \end{aligned}$$

and the confidence interval is

or for difference of 2 means

calculated from all populations

$$\bar{y}_1 - \bar{y}_2 \pm t_{\alpha/2} s \sqrt{\frac{1}{n_1} + \frac{1}{n_2}}.$$

You may have seen this formula before. The linear models framework has allowed us to derive it from first principles.

Exam marks example

We find a confidence interval for the estimable quantity $\mu + \tau_1$, the mean mark of class 1.

```
> tt <- as.vector(c(1,1,0,0))
> halfwidth <- qt(0.975,df=n-k)*sqrt(s2*t(tt) %*% X_t X_c %*% tt)
> tt %*% b + c(-1,1)*halfwidth
[1] 75.68796 84.11204
```

do with lm function

↓ 1 as a factor

```
> newdata <- data.frame(class.f=factor(1))
> predict(model, newdata, interval="confidence", level=0.95)
```

	fit	lwr	upr
1	79.9	75.68796	84.11204

Exam marks example

We find a *prediction* interval for a new student from class 1.

```
> tt <- as.vector(c(1,1,0,0))
> halfwidth <- qt(0.975,df=n-k)*sqrt(s2) *
+      sqrt(1+t(tt)%*%XtXc%*%tt)
> tt %*% b + c(-1,1)*halfwidth

[1] 65.93024 93.86976

> newdata <- data.frame(class.f=factor(1))
> predict(model, newdata, interval="prediction", level=0.95)

    fit      lwr      upr
1 79.9 65.93024 93.86976
```

Exam marks example

We now find a confidence interval for the estimable quantity $\tau_1 - \tau_2$, the difference between the first two classes.

```
> tt <- as.vector(c(0,1,-1,0))
> halfwidth <- qt(0.975,df=n-k)*sqrt(s2*t(tt)%%XtXc%*%tt)
> tt %*% b + c(-1,1)*halfwidth
[1] -12.5567252 -0.6432748
```

```
> confint(model, level=0.95)
```

		2.5 %	97.5 %
$\mu + \tau_1$	(Intercept)	75.6879592	84.11204
$\tau_2 - \tau_1$	class.f2	0.6432748	12.55673
$\tau_3 - \tau_1$	class.f3	3.5432748	15.45673

$$\tau_2 - \tau_1$$

$$\tau_1 - \tau_2 = -(\tau_2 - \tau_1)$$

Exam marks example

We have to express more obscure parameter combinations relative to the treatment contrasts used. Remember:

Label	contr.treatment	contr.sum
Intercept	$\mu + \tau_1$	$\mu + \frac{1}{3} \sum \tau_i$
class.f1		$\tau_1 - \frac{1}{3} \sum \tau_i$
class.f2	$\tau_2 - \tau_1$	$\tau_2 - \frac{1}{3} \sum \tau_i$
class.f3	$\tau_3 - \tau_1$	

So for contr.treatment, $\tau_1 - \tau_2 = -\text{class.f2}$.

```
> library(gmodels)
> ci <- estimable(model, c(0, -1, 0), conf.int=0.95)
> c(ci$Lower, ci$Upper)
[1] -12.5567252 -0.6432748
```

for r parameter = intercept, class.f2, class.f3

Exam marks example

For the contr.sum model, we have

$$\text{Intercept} = \mu + \frac{1}{3}(\tau_1 + \tau_2 + \tau_3)$$

$$\text{class.f1} = \frac{2}{3}\tau_1 - \frac{1}{3}(\tau_2 + \tau_3)$$

$$\text{class.f2} = \frac{2}{3}\tau_2 - \frac{1}{3}(\tau_1 + \tau_3)$$

$$\tau_1 - \tau_2 = \text{class.f1} - \text{class.f2}$$

```
> ci2 <- estimable(model2, c(0,1,-1), conf.int=0.95)
> c(ci2$Lower, ci2$Upper)
[1] -12.5567252 -0.6432748
```

Exam marks example

To find the difference between class 3 and the average of the other two classes, we need

$$\begin{aligned}\tau_3 - \frac{1}{2}\tau_2 - \frac{1}{2}\tau_1 &= (\tau_3 - \tau_1) - \frac{1}{2}(\tau_2 - \tau_1) \\ &= \text{class.f3} - \frac{1}{2}\text{class.f2}.\end{aligned}$$

```
> ci3 <- estimable(model, c(0,-0.5,1), conf.int=0.95)
> c(ci3$Lower, ci3$Upper)
[1] 1.041325 11.358675
```

Exam marks example

For `contr.sum`:

$$\begin{aligned} \text{class.f1} + \text{class.f2} &= \frac{1}{3}\tau_1 + \frac{1}{3}\tau_2 - \frac{2}{3}\tau_3 \\ \tau_3 - \frac{1}{2}\tau_2 - \frac{1}{2}\tau_1 &= -\frac{3}{2}\text{class.f1} - \frac{3}{2}\text{class.f2} \end{aligned}$$

```
> ci4 <- estimable(model2, c(0,-1.5,-1.5), conf.int=0.95)
> c(ci4$Lower, ci4$Upper)
[1] 1.041325 11.358675
```