

Moving between the two Linear Parameterisations in R

Dr.Arnab Chakraborty *

Mrinmoy Banik †

31th August 2022

Abstract

So far we have seen how the same linear estimation problem of finding the yields of various crop varieties can be formulated in 2 different ways. We had also fitted 2 linear models with and without an intercept term. In this section we will see how to switch from the 'mean-yield-type' parameterisations to any of the 'benchmark-yield-type' ones.

Refer to the previous video:45 Keywords: parameterisation, estimation, linear

1 Recap

In the following agri dataset we had three different types of HYV seeds and the corresponding yields in each case. The general linear model is given by:-

$$Y_{ij} = \alpha_i + \epsilon_{ij} \quad (1)$$

But many a times we use a different formulation by introducing a mean yield intercept term μ :-

$$Y_{ij} = \mu + \alpha_i + \epsilon_{ij} \quad (2)$$

Naturally μ then stands for the avg crop yield in any of the varieties with α_i 's indicating the extra yield due to that particular crop variety.

But then we run into the problem of *identifiability*, which can be easily handled by imposing an extra condition on α_i 's like:

$$\sum_{i=1}^3 \alpha_i = 0 \quad \text{or} \quad \alpha_1 = 0$$

TABLE 1: A small dataset of agricultural outputs in 8 separate plots of land.

Agri Dataset:-

SL.no	Variety:-	Yield:-
1	1	210.3
2	2	245
3	2	248.9
4	3	212.3
5	3	230.4
6	2	250.1
7	1	213.5
8	1	212.4

*Department of Applied Statistical Division, Indian Statistical Institute, Kolkata Email:<http://www.isical.ac.in/~arn-abc/>

†This portion of the book is latex-ed and edited by him, Student of B.Stat 3rd yr at ISI, BS2012

2 R Code and Explanation

First we need to set our working directory correctly and import the `agri.csv` dataset. Then fit a linear model, `fit1` using the `lm()` function. In the 1st case we use the formula $yield \sim variety$ which lets R compute the intercept term (μ) by default.

But since the design matrix is only of rank 3, R drops the 1st variety column while fitting the model as we can see from its output below in FIGURE 2 (this is equivalent to setting $\alpha_1 = 0$ in equation (2)).

In `fit2` on the other hand,

$$yield \sim variety - 1$$

is fitted forcing R to drop the intercept term.

By looking at the design matrices we can see they are both linearly independent in both cases.

```

1 #import the dataset
2 agri=read.csv('Agri.csv',header=TRUE)
3 #convert variety column into factor
4 agri$variety=factor(agri$variety)
5 #print the dataset
6 agri
7 #show the dimensions
8 dim(agri)
9 #print summary
10 summary(agri)
11 #plot the box plots
12 plot(agri)
13 #fit the 1st model
14 fit1=lm(yield~variety , agri)
15 #print the estimated coefficients
16 fit1
17 #print the model matrix
18 model.matrix(fit1)
19 #fit the 2nd model
20 fit2=lm(yield~variety-1, agri)
21 #print the model coefficients
22 fit2
23 #print the model matrix
24 model.matrix(fit2)

```

FIGURE 1: The dataset in csv format is imported and two models `fit1` & `fit2` are fitted with `lm()` in R.

Now observe that the intercept in `fit1` and the mean yield or coefficient of `variety1` in `fit2` are essentially same. (the small difference is due to numerical approximation.) This is because `variety1` is the benchmark crop in `fit1` so the intercept μ in this case equal to the mean yield of `variety1` and the rest two coefficients are excess yield of varieties 2 & 3 when compared to that of `variety1`.

Hence the two models, `fit1` & `fit2` are equivalent or bijective parametrization. So the mean yield of `variety2` is 35.993 units greater than the benchmark while `variety3`'s yield is only 9.283 units more making `variety2` the highest yielding crop variety followed by `variety2` and `variety1` is the least productive.

Lastly we can easily switch from one model to the other by subtracting the mean yield of any of the crop variety we want to set as benchmark from the rest and setting its own $\alpha_i = 0$.

```

1 > #fit the 1st model
2 > fit1=lm(yield~variety , agri)
3 > #print the estimated coefficients
4 > fit1
5
6 Call:
7 lm(formula = yield ~ variety , data = agri)
8
9 Coefficients:
10 (Intercept)    variety1    variety3
11 248.00      -35.93      -26.65
12
13 > #print the model matrix
14 > model.matrix(fit1)
15 (Intercept) variety1 variety3
16 1          1          1          0
17 2          1          0          0
18 3          1          0          0
19 4          1          0          1
20 5          1          0          1
21 6          1          0          0
22 7          1          1          0
23 8          1          1          0

```

(A) variety1 is dropped by R to match the rank of the design matrix

```

1 > #fit the 2nd model
2 > fit2=lm(yield~variety -1, agri)
3 > #print the model coefficients
4 > fit2
5
6 Call:
7 lm(formula = yield ~ variety - 1, data =
8   agri)
9
10 Coefficients:
11    variety2    variety1    variety3
12 248.0      212.1      221.3
13
14 > #print the model matrix
15 > model.matrix(fit2)
16    variety2    variety1    variety3
17 1          0          1          0
18 2          1          0          0
19 3          1          0          0
20 4          0          0          1
21 5          0          0          1
22 6          1          0          0
23 7          0          1          0
24 8          0          1          0

```

(B) With no intercept term, all three varieties are fit

FIGURE 2: Fitting the two parametrisations with & without the intercept

3 Interpretation

The estimated coefficients of α_i 's in fit2 are essentially the mean yield of each variety. While in case of fit1, the $\alpha_1 = 0$ so the intercept term basically signifies the benchmark yield of crop 1 which is same as the estimate of the coefficient of variety1 in fit2.

The other two coefficients of variety2 and variety3 in fit1 is equal to $\alpha_2 - \alpha_1$ & $\alpha_3 - \alpha_1$ respectively from fit2. Similarly if we want R to drop variety 2 instead we can use the *relevel()* command as shown in FIGURE 3.

```

1 #relevel variety around variety2
2 agri$variety=relevel(agri$variety , ref=2)
3 #fit the 3rd model
4 fit3=lm(yield~variety , agri)
5 #print the estimated coefficients
6 fit3
7 #print the model matrix
8 model.matrix(fit3)

```

```

1 > model.matrix(fit3)
2 (Intercept)    variety1    variety3
3 1          1          1          0
4 2          1          0          0
5 3          1          0          0
6 4          1          0          1
7 5          1          0          1
8 6          1          0          0
9 7          1          1          0
10 8          1          1          0

```

FIGURE 3: *relevel(agri\$variety , ref=2)* forces R to centre around variety2, hence now it drops variety2 equivalently setting $\alpha_2 = 0$

4 Conclusion

We see there are several ways to interpret the same (mathematically and statistically) linear model. Which one of them is best depends on the situation, but it doesn't change the final inference like here variety2 was superior in yield followed by variety3 and variety1.

5 Appendix

The complete R code for this section is provided below for reference:-

```

1 library(faraway)
2 setwd("C:/My Data/B. Stat/Semester 5/Linear Models/Assignments/Codes")
3 #####video 46#####
4 #import the dataset
5 agri=read.csv('Agri.csv',header=TRUE)
6 #convert variety column into factor
7 agri$variety=factor(agri$variety)
8 #print the dataset
9 agri
10 #show the dimensions
11 dim(agri)
12 #print summary
13 summary(agri)
14 #plot the box plots
15 plot(agri)
16 #fit the 1st model
17 fit1=lm(yield~variety , agri)
18 #print the estimated coefficients
19 fit1
20 #print the model matrix
21 model.matrix(fit1)
22 #fit the 2nd model
23 fit2=lm(yield~variety -1, agri)
24 #print the model coefficients
25 fit2
26 #print the model matrix
27 model.matrix(fit2)
28 #relevel variety around variety2
29 agri$variety=relevel(agri$variety , ref=2)
30 #fit the 3rd model
31 fit3=lm(yield~variety , agri)
32 #print the estimated coefficients
33 fit3
34 #print the model matrix
35 model.matrix(fit3)

```

[—————]