

# Multicollinearity, identification, and estimable functions

Simen Gaure

ABSTRACT. Since there is quite a lot of confusion here and there about what happens when factors are collinear; here is a walkthrough of the identification problems which may arise in models with many dummies, and how **lfe** handles them. (Or, at the very least, attempts to handle them).

## 1. Context

The **lfe** package is used for ordinary least squares estimation, i.e. models which conceptually may be estimated by **lm** as

```
lm(y ~ x1 + x2 + ... + xm + f1 + f2 + ... + fn)
```

where **f1,f2,...,fn** are factors. The standard method is to introduce a dummy variable for each level of each factor. This is too much as it introduces multicollinearities in the system. Conceptually, the system may still be solved, but there are many different solutions. In all of them, the difference between the coefficients for each factor will be the same.

The ambiguity is typically solved by removing a single dummy variable for each factor, this is termed a reference. This is like forcing the coefficient for this dummy variable to zero, and the other levels are then seen as relative to this zero. Other ways to solve the problem is to force the sum of the coefficients to be zero, or one may enforce some other constraint, typically via the **contrasts** argument to **lm**. The default in **lm** is to have a reference level in each factor, and a common intercept term.

In **lfe** the same estimation can be performed by

```
felm(y ~ x1 + x2 + ... + xm | f1 + f2 + ... + fn)
```

Since **felm** conceptually does exactly the same as **lm**, the contrasts approach may work there too. Or rather, it is actually not necessary that **felm** handles it at all, it is only necessary if one needs to fetch the coefficients for the factor levels with **getfe**.

**lfe** is intended for very large datasets, with factors with many levels. Then the approach with a single constraint for each factor may sometimes not be sufficient. The standard example in the econometrics literature (see e.g. [2]) is the case with two factors, one for individuals, and one for firms these individuals work for, changing jobs now and then. What happens in practice is that the labour market may be disconnected, so that one set of individuals move between one set of firms, and

another (disjoint) set of individuals move between some other firms. This happens for no obvious reason, and is data dependent, not intrinsic to the model. There may be several such components. I.e. there are more multicollinearities in the system than the obvious ones. In such a case, there is no way to compare coefficients from different connected components, it is not sufficient with a single individual reference. The problem may be phrased in graph theoretic terms (see e.g. [1, 3, 4]), and it can be shown that it is sufficient with one reference level in each of the connected components. This is what **lfe** does, in the case with two factors it identifies these components, and force one level to zero in one of the factors.

In the examples below, rather small randomly generated datasets are used. **lfe** is hardly the best solution for these problems, they are solely used to illustrate some concepts. I can assure the reader that no CPUs, sleeping patterns, romantic relationships, trees or cats, nor animals in general, were harmed during data collection and analysis.

## 2. Identification with two factors

In the case with two factors, identification is well-known. **getfe** will partition the dataset into connected components, and introduce a reference level in each component:

```
library(lfe)
## Loading required package: Matrix
set.seed(42)
x1 <- rnorm(20)
f1 <- sample(8,length(x1),replace=TRUE)/10
f2 <- sample(8,length(x1),replace=TRUE)/10
e1 <- sin(f1) + 0.02*f2^2 + rnorm(length(x1))
y <- 2.5*x1 + (e1-mean(e1))
summary(est <- felm(y ~ x1 | f1 + f2))
##
## Call:
##   felm(formula = y ~ x1 | f1 + f2)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.399 -0.279  0.000  0.436  0.981
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## x1      2.531      0.377   6.71  0.0011 **
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.13 on 5 degrees of freedom
## Multiple R-squared:  0.973    Adjusted R-squared:  0.894
## F-statistic:13.1 on 14 and 5 DF, p-value: 0.0051
```

We examine the estimable function produced by **efactory**.

```

ef <- efactory(est)
is.estimable(ef,est$fe)
## [1] TRUE
getfe(est)
##          effect obs comp fe idx
## f1.0.1  0.376275  2    1 f1 0.1
## f1.0.2 -0.081100  1    2 f1 0.2
## f1.0.3 -0.686880  3    1 f1 0.3
## f1.0.4  0.573177  4    1 f1 0.4
## f1.0.5  0.479142  2    1 f1 0.5
## f1.0.6  1.413020  3    1 f1 0.6
## f1.0.7  0.844956  1    2 f1 0.7
## f1.0.8  0.926434  4    1 f1 0.8
## f2.0.1 -0.004011  3    1 f2 0.1
## f2.0.2  0.000000  5    1 f2 0.2
## f2.0.3 -1.518667  1    1 f2 0.3
## f2.0.4  0.000000  2    2 f2 0.4
## f2.0.5 -1.894524  2    1 f2 0.5
## f2.0.6 -0.884319  3    1 f2 0.6
## f2.0.7 -0.609110  3    1 f2 0.7
## f2.0.8 -0.968652  1    1 f2 0.8

```

As we can see from the `comp` entry, there are two components, with `f1=0.2`, `f1=0.7` and `f2=0.4`. A reference is introduced in each of the components, i.e. `f2.0.2=0` and `f2.0.4=0`. If we look at the dataset, the component structure becomes clearer:

```

data.frame(f1,f2,comp=est$cfactor)
##      f1  f2 comp
## 1  0.4 0.6   1
## 2  0.4 0.8   1
## 3  0.1 0.7   1
## 4  0.8 0.5   1
## 5  0.4 0.7   1
## 6  0.8 0.2   1
## 7  0.8 0.3   1
## 8  0.6 0.7   1
## 9  0.8 0.6   1
## 10 0.5 0.2   1
## 11 0.3 0.1   1
## 12 0.3 0.2   1
## 13 0.4 0.2   1
## 14 0.7 0.4   2
## 15 0.1 0.2   1
## 16 0.6 0.6   1
## 17 0.6 0.1   1
## 18 0.2 0.4   2
## 19 0.3 0.5   1

```

```
## 20 0.5 0.1 1
```

Observation 14 and 18 belong to component 2; no other observation has  $f1=0.7$ ,  $f1=0.2$  or  $f2=0.4$ , thus it is clear that coefficients for these can not be compared to other coefficients. `lm` is silent about this component structure, hence coefficients are hard to interpret. Though, predictive properties and residuals are the same:

```
f1 <- factor(f1); f2 <- factor(f2)
summary(lm(y ~ x1 + f1 + f2))
##
## Call:
## lm(formula = y ~ x1 + f1 + f2)
##
## Residuals:
```

	1	2	3	4	5	6	7
##	4.18e-01	4.72e-16	6.36e-01	4.91e-01	-1.40e+00	-2.26e-01	1.94e-16
##	8	9	10	11	12	13	14
##	7.64e-01	-2.64e-01	2.05e-01	8.15e-01	-3.24e-01	9.81e-01	1.53e-16
##	15	16	17	18	19	20	
##	-6.36e-01	-1.54e-01	-6.10e-01	-1.25e-16	-4.91e-01	-2.05e-01	

```
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  0.37226    1.21296   0.31   0.7713
## x1           2.53052    0.37710   6.71   0.0011 **
## f10.2        -0.45336    2.02866  -0.22   0.8320
## f10.3        -1.06316    1.24867  -0.85   0.4334
## f10.4         0.19690    1.07494   0.18   0.8619
## f10.5         0.10287    1.27009   0.08   0.9386
## f10.6         1.03674    1.15377   0.90   0.4101
## f10.7         0.47269    1.67007   0.28   0.7885
## f10.8         0.55016    1.25945   0.44   0.6805
## f20.2         0.00401    0.96913   0.00   0.9969
## f20.3        -1.51466    1.56345  -0.97   0.3771
## f20.4         NA         NA         NA      NA
## f20.5        -1.89051    1.50465  -1.26   0.2645
## f20.6        -0.88031    1.11170  -0.79   0.4643
## f20.7        -0.60510    1.09728  -0.55   0.6051
## f20.8        -0.96464    1.58571  -0.61   0.5695
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.13 on 5 degrees of freedom
## Multiple R-squared:  0.973, Adjusted R-squared:  0.899
## F-statistic: 13.1 on 14 and 5 DF, p-value: 0.0051
```

### 3. Identification with three or more factors

In the case with three or more factors, there is no general intuitive theory (yet) for handling identification problems. **lfe** resorts to the simple-minded approach that non-obvious multicollinearities arise among the first two factors, and assumes it is sufficient with a single reference level for each of the remaining factors, i.e. that they in principle could be specified as ordinary dummies. In other words, the order of the factors in the model specification is important. A typical example would be 3 factors; individuals, firms and education:

```
est <- felm(logwage ~ x1 + x2 | id + firm + edu)
getfe(est)
```

This will result in the same number of references as if using the model

```
logwage ~ x1 + x2 + edu | id + firm
```

though it may run faster (or slower).

Alternatively, one could specify the model as

```
logwage ~ x1 + x2 | firm + edu + id
```

This would not account for a partitioning of the labour market along individual/firm, but along firm/education, using a single reference level for the individuals. In this example, there is some reason to suspect that it is not sufficient, depending on how **edu** is specified. There exists no general scheme that sets up suitable reference groups when there are more than two factors. It may happen that the default is sufficient. The function **getfe** will check whether this is so, and it will yield a warning about 'non-estimable function' if not. With some luck it may be possible to rearrange the order of the factors to avoid this situation.

There is nothing special with **lfe** in this respect. You will meet the same problem with **lm**, it will remove a reference level (or dummy-variable) in each factor, but the system will still contain multicollinearities. You may remove reference levels until all the multicollinearities are gone, but there is no obvious way to interpret the resulting coefficients.

To illustrate, the classical example is when you include a factor for age (in years), a factor for observation year, and a factor for year of birth. You pick a reference individual, e.g. **age=50**, **year=2013** and **birth=1963**, but this is not sufficient to remove all the multicollinearities. If you analyze this problem (see e.g. [6]) you will find that the coefficients are only identified up to linear trends. You may force the linear trend between **birth=1963** and **birth=1990** to zero, by removing the reference level **birth=1990**, and the system will be free of multicollinearities. In this case the **birth** coefficients have the interpretation as being deviations from a linear trend between 1963 and 1990, though you do not know which linear trend. The **age** and **year** coefficients are also relative to this same unknown trend.

In the above case, the multicollinearity is obviously built into the model, and it is possible to remove it and find some intuitive interpretation of the coefficients. In the general case, when either **lm** or **getfe** reports a handful of non-obvious spurious multicollinearities between factors with many levels, you probably will not be able to find any reasonable way to interpret coefficients. Of course, certain linear

combinations of coefficients will be unique, i.e. estimable, and these may be found by e.g. the procedures in [5, 8], but the general picture is muddy.

**lfe** does not provide a solution to this problem, however, **getfe** will still provide a vector of coefficients which results from finding a non-unique solution to a certain set of equations. To get any sense from this, an estimable function must be applied. The simplest one is to pick a reference for each factor and subtract this coefficient from each of the other coefficients in the same factor, and add it to a common intercept, however in the case this does not result in an estimable function, you are out of luck. If you for some reason believe that you know of an estimable function, you may provide this to **getfe** via the **ef**-argument. There is an example in the **getfe** documentation. You may also test it for estimability with the function **is.estimable**, this is a probabilistic test which almost never fails (see [4, Remark 6.2]).

#### 4. Specifying an estimable function

A model of the type

```
y ~ x1 + x2 + f1 + f2 + f3
```

may be written in matrix notation as

$$(1) \quad y = X\beta + D\alpha + \epsilon,$$

where  $X$  is a matrix with columns **x1** and **x2** and  $D$  is matrix of dummies constructed from the levels of the factors **f1**, **f2**, **f3**. Formally, an estimable function in our context is a matrix operator whose row space is contained in the row space of  $D$ . That is, an estimable function may be written as a matrix. Like the **contrasts** argument to **lm**. However, the **lfe** package uses an R-function instead. That is, **felm** is called first, it uses the Frisch-Waugh-Lovell theorem to project out the  $D\alpha$  term from (1) (see [4, Remark 3.2]):

```
est <- felm(y ~ x1 + x2 | f1 + f2 + f3)
```

This yields the parameters for **x1** and **x2**, i.e.  $\hat{\beta}$ . To find  $\hat{\alpha}$ , the parameters for the levels of **f1**, **f2**, **f3**, **getfe** solves a certain linear system (see [4, eq. (14)]):

$$(2) \quad D\gamma = \rho$$

where the vector  $\rho$  can be computed when we have  $\hat{\beta}$ . This does not identify  $\gamma$  uniquely, we have to apply an estimable function to  $\gamma$ . The estimable function  $F$  is characterized by the property that  $F\gamma_1 = F\gamma_2$  whenever  $\gamma_1$  and  $\gamma_2$  are solutions to equation (2). Rather than coding  $F$  as a matrix, **lfe** codes it as a function. It is of course possible to let the function apply a matrix, so this is not a material distinction. So, let's look at an example of how an estimable function may be made:

```
library(lfe)
x1 <- rnorm(100)
f1 <- sample(7,100,replace=TRUE)
f2 <- sample(8,100,replace=TRUE)/8
f3 <- sample(10,100,replace=TRUE)/10
e1 <- sin(f1) + 0.02*f2^2 + 0.17*f3^3 + rnorm(100)
y <- 2.5*x1 + (e1-mean(e1))
```

```
summary(est <- felm(y ~ x1 | f1 + f2 + f3))
##
## Call:
##   felm(formula = y ~ x1 | f1 + f2 + f3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.4435 -0.7097  0.0205  0.7724  1.9945
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## x1      2.453      0.114    21.6   <2e-16 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.13 on 76 degrees of freedom
## Multiple R-squared:  0.881    Adjusted R-squared:  0.844
## F-statistic:24.6 on 23 and 76 DF, p-value: <2e-16
## *** Standard errors may be too high due to more than 2 groups and exactDOF=FALSE
```

In this case, with 3 factors we can not be certain that it is sufficient with a single reference in two of the factors, but we try it as an exercise. (**lfe** does not include an intercept, it is subsumed in one of the factors, so it should tentatively be sufficient with a reference for the two others).

The input to our estimable function is a solution  $\gamma$  of equation (2). The argument **addnames** is a logical, set to **TRUE** when the function should add names to the resulting vector. The coefficients is ordered the same way as the levels in the factors. We should pick a single reference in factors **f2**, **f3**, subtract these, and add the sum to the first factor:

```
ef <- function(gamma, addnames) {
  ref2 <- gamma[[8]]
  ref3 <- gamma[[16]]
  gamma[1:7] <- gamma[1:7] + ref2 + ref3
  gamma[8:15] <- gamma[8:15] - ref2
  gamma[16:25] <- gamma[16:25] - ref3
  if(addnames) {
    names(gamma) <- c(paste('f1', 1:7, sep='.'),
                      paste('f2', 1:8, sep='.'),
                      paste('f3', 1:10, sep='.'))
  }
  gamma
}
is.estimable(ef, fe=est$fe)
## [1] TRUE
getfe(est, ef=ef)
##      effect
## f1.1 -0.23664
```

```
## f1.2 -0.03447
## f1.3 -0.90006
## f1.4 -1.88954
## f1.5 -2.47066
## f1.6 -2.02413
## f1.7 -0.78115
## f2.1 0.00000
## f2.2 0.20650
## f2.3 0.05873
## f2.4 1.14101
## f2.5 0.19097
## f2.6 0.66239
## f2.7 0.25344
## f2.8 0.29348
## f3.1 0.00000
## f3.2 0.82025
## f3.3 0.56127
## f3.4 0.80081
## f3.5 0.81496
## f3.6 0.44563
## f3.7 1.36639
## f3.8 0.46518
## f3.9 1.32877
## f3.10 0.97592
```

We may compare this to the default estimable function, which picks a reference in each connected component as defined by the two first factors.

```
getfe(est)
##          effect obs comp fe    idx
## f1.1      0.00000 17   1 f1     1
## f1.2      0.20218 15   1 f1     2
## f1.3     -0.66342 16   1 f1     3
## f1.4     -1.65290 15   1 f1     4
## f1.5     -2.23402 12   1 f1     5
## f1.6     -1.78749 12   1 f1     6
## f1.7     -0.54451 13   1 f1     7
## f2.0.125  0.22854 12   1 f2 0.125
## f2.0.25   0.43504 15   1 f2 0.25
## f2.0.375  0.28728 13   1 f2 0.375
## f2.0.5    1.36956 13   1 f2 0.5
## f2.0.625  0.41951 12   1 f2 0.625
## f2.0.75   0.89094 13   1 f2 0.75
## f2.0.875  0.48198 11   1 f2 0.875
## f2.1       0.52202 11   1 f2     1
## f3.0.1   -0.46518  4   2 f3 0.1
## f3.0.2    0.35506 10   2 f3 0.2
## f3.0.3    0.09609 14   2 f3 0.3
## f3.0.4    0.33562 13   2 f3 0.4
```



```
## f3.0.5    0.34978    9    2 f3    0.5
## f3.0.6   -0.01955   10    2 f3    0.6
## f3.0.7    0.90120    6    2 f3    0.7
## f3.0.8    0.00000   16    2 f3    0.8
## f3.0.9    0.86358   11    2 f3    0.9
## f3.1      0.51073    7    2 f3     1
```

We see that the default has some more information. It uses the level names, and some more information, added like this:

```
efactory(est)
## function (v, addnames)
## {
##   esum <- sum(v[extrarefs])
##   df <- v[refsubs]
##   sub <- ifelse(is.na(df), 0, df)
##   df <- v[refsuba]
##   add <- ifelse(is.na(df), 0, df + esum)
##   v <- v - sub + add
##   if (addnames) {
##     names(v) <- nm
##     attr(v, "extra") <- list(obs = obs, comp = comp, fe = fef,
##                               idx = idx)
##   }
##   v
## }
## <bytecode: 0x743dd80>
## <environment: 0x770c3a0>
```

I.e. when asked to provide level names, it is also possible to add additional information as a list (or `data.frame`) as an attribute `'extra'`. The vectors `extrarefs`, `refsubs`, `refsuba` etc. are precomputed by `efactory` for speed efficiency.

Here is the above example, but we create an intercept instead, and don't report the zero-coefficients, so that it closely resembles the output from `lm`

```
f1 <- factor(f1); f2 <- factor(f2); f3 <- factor(f3)
ef <- function(gamma, addnames) {
  ref1 <- gamma[[1]]
  ref2 <- gamma[[8]]
  ref3 <- gamma[[16]]
  # put the intercept in the first coordinate
  gamma[[1]] <- ref1+ref2+ref3
  gamma[2:7] <- gamma[2:7]-ref1
  gamma[8:14] <- gamma[9:15]-ref2
  gamma[15:23] <- gamma[17:25]-ref3
  length(gamma) <- 23
  if(addnames) {
    names(gamma) <- c('(Intercept)', paste('f1', levels(f1)[2:7], sep=' '),
                     paste('f2', levels(f2)[2:8], sep=' '),
```

```

    paste('f3',levels(f3)[2:10],sep=''))
  }
  gamma
}
getfe(est,ef=ef,bN=1000,se=TRUE)
##              effect      se
## (Intercept) -0.23664 0.7693
## f12          0.20218 0.4590
## f13         -0.66342 0.4223
## f14         -1.65290 0.4858
## f15         -2.23402 0.4560
## f16         -1.78749 0.4625
## f17         -0.54451 0.4882
## f20.25       0.20650 0.4435
## f20.375      0.05873 0.5073
## f20.5        1.14101 0.5067
## f20.625      0.19097 0.5289
## f20.75       0.66239 0.4643
## f20.875      0.25344 0.5199
## f21          0.29348 0.4917
## f30.2        0.82025 0.6765
## f30.3        0.56127 0.6882
## f30.4        0.80081 0.7017
## f30.5        0.81496 0.7288
## f30.6        0.44563 0.7642
## f30.7        1.36639 0.7745
## f30.8        0.46518 0.6815
## f30.9        1.32877 0.7117
## f31          0.97592 0.7888

#compare with lm
summary(lm(y ~ x1 + f1 + f2 + f3))

##
## Call:
## lm(formula = y ~ x1 + f1 + f2 + f3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -2.4435 -0.7097  0.0205  0.7724  1.9945
##
## Coefficients:
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -0.2366      0.7698  -0.31   0.7594
## x1            2.4530      0.1136  21.58 <2e-16 ***
## f12           0.2022      0.4444   0.45   0.6504
## f13          -0.6634      0.4312  -1.54   0.1281
## f14          -1.6529      0.4866  -3.40   0.0011 **
## f15          -2.2340      0.4727  -4.73   1e-05 ***

```

```
## f16      -1.7875      0.4822     -3.71      0.0004 ***
## f17      -0.5445      0.4957     -1.10      0.2754
## f20.25    0.2065      0.4564      0.45      0.6523
## f20.375   0.0587      0.4959      0.12      0.9060
## f20.5     1.1410      0.4926      2.32      0.0232 *
## f20.625   0.1910      0.5122      0.37      0.7103
## f20.75    0.6624      0.4788      1.38      0.1705
## f20.875   0.2534      0.5273      0.48      0.6321
## f21       0.2935      0.5081      0.58      0.5653
## f30.2     0.8202      0.7021      1.17      0.2464
## f30.3     0.5613      0.6881      0.82      0.4172
## f30.4     0.8008      0.7117      1.13      0.2640
## f30.5     0.8150      0.7255      1.12      0.2648
## f30.6     0.4456      0.7438      0.60      0.5509
## f30.7     1.3664      0.7915      1.73      0.0883 .
## f30.8     0.4652      0.6708      0.69      0.4901
## f30.9     1.3288      0.7132      1.86      0.0663 .
## f31       0.9759      0.7830      1.25      0.2164
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.13 on 76 degrees of freedom
## Multiple R-squared:  0.881, Adjusted R-squared:  0.846
## F-statistic: 24.6 on 23 and 76 DF, p-value: <2e-16
```

## 5. Non-estimability

We consider another example. To ensure spurious relations there are almost as many factor levels as there are observations, and it will be hard to find enough estimable function to interpret all the coefficients. The coefficient for  $x_1$  is still estimated, but with a large standard error. Note that this is an illustration of non-obvious non-estimability which may occur in much larger datasets, the author does not endorse using this kind of model for the kind of data you find below.

```
set.seed(128)
x1 <- rnorm(25)
f1 <- sample(9,length(x1),replace=TRUE)
f2 <- sample(8,length(x1),replace=TRUE)
f3 <- sample(8,length(x1),replace=TRUE)
e1 <- sin(f1) + 0.02*f2^2 + 0.17*f3^3 + rnorm(length(x1))
y <- 2.5*x1 + (e1-mean(e1))
summary(est <- felm(y ~ x1 | f1 + f2 + f3))
##
## Call:
## felm(formula = y ~ x1 | f1 + f2 + f3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -0.870 -0.368 0.000 0.148 1.884
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## x1      1.41      0.47      3      0.03 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.3 on 5 degrees of freedom
## Multiple R-squared:  1    Adjusted R-squared: 0.999
## F-statistic: 933 on 19 and 5 DF, p-value: 1.34e-07
## *** Standard errors may be too high due to more than 2 groups and exactDOF=FALSE
```

The default estimable function fails, and the coefficients from `getfe` are not useable. `getfe` yields a warning in this case.

```
ef <- efactory(est)
is.estimable(ef, est$fe)
## Warning: non-estimable function, largest error 0.3 in coordinate 4
## ("f1.5")
## [1] FALSE
```

Indeed, the rank-deficiency is larger than expected. There are more spurious relations between the factors than what can be accounted for by looking at components in the two first factors. In this low-dimensional example we may find the matrix  $D$  of equation (2), and its (column) rank deficiency is larger than 2.

```
f1 <- factor(f1); f2 <- factor(f2); f3 <- factor(f3)
D <- t(do.call('rBind',
              lapply(list(f1, f2, f3), as, Class='sparseMatrix'))))
dim(D)
## [1] 25 21
as.integer(rankMatrix(D))
## [1] 18
```

Alternatively we can use an internal function in `lfe` for finding the rank deficiency directly.

```
lfe::rankDefic(list(f1, f2, f3))
## [1] 3
```

This rank-deficiency also has an impact on the standard errors computed by `felm`. If the rank-deficiency is small relative to the degrees of freedom the standard errors are scaled slightly upwards if we ignore the rank deficiency, but if it is large, the impact on the standard errors can be substantial. The above mentioned rank-computation procedure can be activated by specifying `exactDOF=TRUE` in the call to `felm`, but it may be time-consuming if the factors have many levels. Computing the rank does not in itself help us find estimable functions for `getfe`.

```
summary(est <- felm(y ~ x1 | f1 + f2 + f3, exactDOF=TRUE))
```

```
##
## Call:
##   felm(formula = y ~ x1 | f1 + f2 + f3, exactDOF = TRUE)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.870 -0.368  0.000  0.148  1.884
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## x1      1.408      0.429   3.28  0.017 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.19 on 6 degrees of freedom
## Multiple R-squared:  1    Adjusted R-squared: 0.999
## F-statistic:1.18e+03 on 18 and 6 DF, p-value: 3.7e-09
```

We can get an idea what happens if we keep the dummies for **f3**. In this case, with 2 factors, **lfe** will partition the dataset into connected components and account for all the multicollinearities among the factors **f1** and **f2** just as above, but this is not sufficient. The interpretation of the resulting coefficients is not straightforward.

```
summary(est <- felm(y ~ x1 + f3 | f1 + f2, exactDOF=TRUE))
##
## Call:
##   felm(formula = y ~ x1 + f3 | f1 + f2, exactDOF = TRUE)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -8.70e-01 -3.68e-01 -4.42e-08  1.48e-01  1.88e+00
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## x1      1.408      0.429   3.28  0.017 *
## f32      0.697      1.058   0.66  0.535
## f33      4.231      1.544   2.74  0.034 *
## f34         NA         NA     NA     NA
## f35     19.489      1.160  16.81 2.8e-06 ***
## f36     35.684      2.264  15.76 4.1e-06 ***
## f37     57.772      1.185  48.75 5.0e-09 ***
## f38     86.498      1.226  70.53 5.5e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.19 on 6 degrees of freedom
## Multiple R-squared:  1    Adjusted R-squared: 0.999
## F-statistic:1.18e+03 on 18 and 6 DF, p-value: 3.7e-09
```

```
getfe(est)
##          effect obs comp fe idx
## f1.1 -34.5246   4    1 f1   1
## f1.2 -33.3982   5    1 f1   2
## f1.3 -34.6698   6    1 f1   3
## f1.5 -27.3871   1    1 f1   5
## f1.6 -35.7423   5    1 f1   6
## f1.7 -36.4468   1    1 f1   7
## f1.9 -32.2795   3    1 f1   9
## f2.2  -0.4037   3    1 f2   2
## f2.3   1.8848   3    1 f2   3
## f2.5  -2.6807   2    1 f2   5
## f2.6   0.0000   8    1 f2   6
## f2.7   0.8925   4    1 f2   7
## f2.8   0.9553   5    1 f2   8
```

In this particular example, we may use a different order of the factors, and we see that by partitioning the dataset on the factors `f1,f3` instead of `f1,f2`, there are 2 connected components (the factor `f2` gets its own `comp`-code, but this is not a graph theoretic component number, it merely indicates that there is a separate reference among these).

```
summary(est <- felm(y ~ x1 | f1 + f3 + f2, exactDOF=TRUE))
##
## Call:
## felm(formula = y ~ x1 | f1 + f3 + f2, exactDOF = TRUE)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.870 -0.368   0.000   0.148   1.884
##
## Coefficients:
##      Estimate Std. Error t value Pr(>|t|)
## x1      1.408      0.429    3.28   0.017 *
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.19 on 6 degrees of freedom
## Multiple R-squared:  1    Adjusted R-squared:  0.999
## F-statistic:1.18e+03 on 18 and 6 DF, p-value: 3.7e-09
is.estimable(efactory(est),est$fe)
## [1] TRUE
getfe(est)
##          effect obs comp fe idx
## f1.1   0.1452   4    1 f1   1
## f1.2   1.2716   5    1 f1   2
## f1.3   0.0000   6    1 f1   3
## f1.5   0.0000   1    2 f1   5
```

```
## f1.6 -1.0725 5 1 f1 6
## f1.7 -1.7770 1 1 f1 7
## f1.9 2.3903 3 1 f1 9
## f3.1 -34.6698 4 1 f3 1
## f3.2 -33.9730 4 1 f3 2
## f3.3 -30.4384 1 1 f3 3
## f3.4 -27.3871 1 2 f3 4
## f3.5 -15.1809 4 1 f3 5
## f3.6 1.0138 2 1 f3 6
## f3.7 23.1023 3 1 f3 7
## f3.8 51.8278 6 1 f3 8
## f2.2 -0.4037 3 3 f2 2
## f2.3 1.8848 3 3 f2 3
## f2.5 -2.6807 2 3 f2 5
## f2.6 0.0000 8 3 f2 6
## f2.7 0.8925 4 3 f2 7
## f2.8 0.9553 5 3 f2 8
```

Below is the same estimation in `lm`. We see that the coefficient for `x1` is identical to the one from `felm`, but there is no obvious relation between e.g. the coefficients for `f1`; the difference `f15-f16` is not the same for `lm` and `felm`. Since these are in different components, they are not comparable. But of course, if we compare in the same component, e.g. `f16-f17` or take a combination which actually occurs in the dataset, it is unique (estimable):

```
data.frame(f1,f2,f3)[1,]
##   f1 f2 f3
## 1  1  6  7
```

I.e. if we add the coefficients `f1.1 + f2.6 + f3.7` and include the intercept for `lm`, we will get the same number for both `lm` and `felm`.

```
summary(est <- lm(y ~ x1 + f1 + f2 + f3))
##
## Call:
## lm(formula = y ~ x1 + f1 + f2 + f3)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -0.870 -0.368  0.000  0.148  1.884
##
## Coefficients: (1 not defined because of singularities)
##              Estimate Std. Error t value Pr(>|t|)
## (Intercept)  -34.928     0.981  -35.59  3.3e-08 ***
## x1              1.408     0.429   3.28  0.0168 *
## f12             1.126     1.354   0.83  0.4373
## f13            -0.145     1.162  -0.12  0.9047
## f15             7.138     1.822   3.92  0.0078 **
## f16            -1.218     1.412  -0.86  0.4215
## f17            -1.922     2.240  -0.86  0.4237
```

```
## f19      2.245      1.925      1.17      0.2877
## f23      2.289      1.562      1.47      0.1932
## f25     -2.277      1.455     -1.56      0.1687
## f26      0.404      1.187      0.34      0.7454
## f27      1.296      1.528      0.85      0.4289
## f28      1.359      1.477      0.92      0.3930
## f32      0.697      1.058      0.66      0.5346
## f33      4.231      1.544      2.74      0.0337 *
## f34      NA         NA         NA         NA
## f35     19.489      1.160     16.81     2.8e-06 ***
## f36     35.684      2.264     15.76     4.1e-06 ***
## f37     57.772      1.185     48.75     5.0e-09 ***
## f38     86.498      1.226     70.53     5.5e-10 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.19 on 6 degrees of freedom
## Multiple R-squared:      1, Adjusted R-squared:  0.999
## F-statistic: 1.18e+03 on 18 and 6 DF,  p-value: 3.7e-09
```

## 6. Weeks-Williams partitions

There is a partial solution to the non-estimability problem in [8]. Their idea is to partition the dataset into components in which all differences between factor levels are estimable. The components are connected components of a subgraph of an  $e$ -dimensional grid graph where  $e$  is the number of factors. That is, a graph is constructed with the observations as vertices, two observations are adjacent (in a graph theoretic sense) if they differ in at most one of the factors. The dataset is then partitioned into (graph theoretic) connected components. It's a finer partitioning than the above, and consequently introduces more reference levels than is necessary for identification. I.e. it does not find all estimable functions, but in some cases (e.g. in [7]) the largest component will be sufficiently large for proper analysis. It is of course always a question whether such an endogeneous selection of observations will yield a dataset which results in unbiased coefficients. This partitioning can be done by the `compfactor` function with argument `WW=TRUE`:

```
fe <- list(f1,f2,f3)
wwcomp <- compfactor(fe, WW=TRUE)
```

It has more levels than the rank deficiency

```
lfe::rankDefic(fe)
## [1] 3
nlevels(wwcomp)
## [1] 10
```

and each of its components are contained in a component of the previously considered components, no matter which two factors we consider. For the case of two factors, the concepts coincide.



```
nlevels(interaction(compfactor(fe),wwcomp))
## [1] 10
# pick the largest component:
wwdata <- data.frame(y, x1, f1, f2, f3)[wwcomp==1, ]
print(wwdata)
##           y           x1 f1 f2 f3
##  2    52.91  0.482612  2  6  8
##  8    54.99  0.004553  2  6  8
## 11   -13.28  0.322709  2  8  5
## 15    54.62  0.684795  2  7  8
## 16   -35.09 -0.883670  3  6  2
## 17    56.13  0.722127  9  7  8
## 19   -14.60 -0.563255  2  6  5
## 20   -18.94 -1.644455  6  6  5
## 21   -34.37  0.359966  3  2  2
## 25    52.13  0.589060  3  6  8
```

Though, in this particular example, there are more parameters than there are observations, so an estimation would not be feasible.

`efactory` cannot easily be modified to produce an estimable function corresponding to WW components. The reason is that `efactory`, and the logic in `getfe`, work on partitions of factor levels, not on partitions of the dataset, these are the same for the two-factor case.

WW partitions have the property that if you pick any two of the factors and partition a WW-component into the previously mentioned non-WW partitions, there will be only one component, hence you may use any of the estimable functions from `efactory` on each partition. That is, a way to use WW partitions with `lfe` is to do the whole analysis on the largest WW-component. `felm` may still be used on the whole dataset, and it may yield different results than what you get by analysing the largest WW-component.

Here is a larger example:

```
set.seed(135)
x <- rnorm(10000)
f1 <- sample(1000,length(x),replace=TRUE)
f2 <- (f1 + sample(18,length(x), replace=TRUE)) %% 500
f3 <- (f2 + sample(9,length(x),replace=TRUE)) %% 500
y <- x + 1e-4*f1 + sin(f2^2) +
  cos(f3)^3 + 0.5*rnorm(length(x))
dataset <- data.frame(y,x,f1,f2,f3)
summary(est <- felm(y ~ x | f1 + f2 + f3,
  data=dataset, exactDOF=TRUE))
##
## Call:
## felm(formula = y ~ x | f1 + f2 + f3, data = dataset, exactDOF = TRUE)
##
## Residuals:
##      Min       1Q   Median       3Q      Max
```

```
## -1.929087 -0.305385 0.000243 0.305398 1.921153
##
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## x 0.99991 0.00565 177 <2e-16 ***
## ---
## Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.505 on 8001 degrees of freedom
## Multiple R-squared: 0.903 Adjusted R-squared: 0.879
## F-statistic:37.5 on 1998 and 8001 DF, p-value: <2e-16
```

We count the number of connected components in `f1,f2`, and see that this is sufficient to ensure estimability

```
nlevels(est$cfactor)
## [1] 1
is.estimable(efactory(est), est$fe)
## [1] TRUE
nrow(alpha <- getfe(est))
## [1] 2000
```

It has rank deficiency one less than the number of factors :

```
lfe::rankDefic(est$fe)
## [1] 2
```

Then we analyse the largest WW-component

```
wwcomp <- compfactor(est$fe,WW=TRUE)
nlevels(wwcomp)
## [1] 958
wwset <- wwcomp == 1
sum(wwset)
## [1] 2394
summary(wwest <- felm(y ~ x | f1 + f2 + f3,
                     data=dataset, subset=wwset, exactDOF=TRUE))
##
## Call:
## felm(formula = y ~ x | f1 + f2 + f3, data = dataset, exactDOF = TRUE, subset = wws
##
## Residuals:
##      Min       1Q   Median       3Q      Max
## -1.882864 -0.294976  0.000194  0.294391  2.037462
##
## Coefficients:
## Estimate Std. Error t value Pr(>|t|)
## x      1.003      0.012   83.5 <2e-16 ***
## ---
```

```
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 0.514 on 1788 degrees of freedom
## Multiple R-squared:  0.911    Adjusted R-squared:  0.881
## F-statistic:30.3 on 605 and 1788 DF, p-value: <2e-16
```

We see that we get the same coefficient for  $x$  in this case. This is not surprising, there is no obvious reason to believe that our selection of observations is skewed in this randomly created dataset.

This one has the same rank deficiency:

```
lfe::rankDefic(wwest$fe)
## [1] 2
```

but a smaller number of identifiable coefficients.

```
nrow(wwalpha <- getfe(wwest))
## [1] 607
```

We may compare effects which are common to the two methods:

```
head(alpha)
##      effect obs comp fe idx
## f1.1 -0.5542  6    1 f1   1
## f1.2 -0.2756 13    1 f1   2
## f1.3 -0.5763  6    1 f1   3
## f1.4 -0.5413  8    1 f1   4
## f1.5 -0.5697 13    1 f1   5
## f1.6 -0.4250  9    1 f1   6
head(wwalpha)
##      effect obs comp fe idx
## f1.1 -1.27887  3    1 f1   1
## f1.2 -0.34865  7    1 f1   2
## f1.3 -0.95120  4    1 f1   3
## f1.4 -1.04719  1    1 f1   4
## f1.6 -0.89773  4    1 f1   6
## f1.7  0.04965  1    1 f1   7
```

but there is no obvious relation between e.g.  $f1.1 - f1.2$ , they are very different in the two estimations. The coefficients are from different datasets, and the standard errors are large ( $\approx 0.7$ ) with this few observations for each factor level. The number of identified coefficients for each factor varies (these figures contain the two references):

```
table(wwalpha[, 'fe'])
##
##  f1  f2  f3
## 310 148 149
```

## References

- [1] J. M. Abowd, R. H. Creecy, and F. Kramarz, *Computing Person and Firm Effects Using Linked Longitudinal Employer-Employee Data*, Tech. Report TP-2002-06, U.S. Census Bureau, 2002.
- [2] J. M. Abowd, F. Kramarz, and D. N. Margolis, *High Wage Workers and High Wage Firms*, *Econometrica* **67** (1999), no. 2, 251–333.
- [3] J. A. Eccleston and A. Hedayat, *On the Theory of Connected Designs: Characterization and Optimality*, *Ann. Statist.* **2** (1974), 1238–1255.
- [4] S. Gaure, *OLS with Multiple High Dimensional Category Variables*, *Computational Statistics and Data Analysis* **66** (2013), 8–18.
- [5] J. D. Godolphin and E. J. Godolphin, *On the connectivity of row-column designs*, *Util. Math.* **60** (2001), 51–65.
- [6] L.L. Kupper, J.M. Janis, I.A. Salama, C.N. Yoshizawa, and B.G. Greenberg, *Age-Period-Cohort Analysis: An Illustration of the Problems in Assessing Interaction in One Observation Per Cell Data*, *Commun. Statist.-Theor. Meth.* **12** (1983), no. 23, 2779–2807.
- [7] S.M. Torres, P. Portugal, J.T. Addison, and P. Guimarães, *The Sources of Wage Variation: A Three-Way High-Dimensional Fixed Effects Regression Model.*, IZA Discussion Paper 7276, Institute for the Study of Labor (IZA), March 2013.
- [8] D.L. Weeks and D.R. Williams, *A Note on the Determination of Connectedness in an N-Way Cross Classification*, *Technometrics* **6** (1964), no. 3, 319–324.

RAGNAR FRISCH CENTRE FOR ECONOMIC RESEARCH, OSLO, NORWAY