# factorMerger: hierarchical clustering and model visualization

Agnieszka Sitko

University of Warsaw, MI^2 Group

useR!2017 | 06-07-2017

## Problem 1

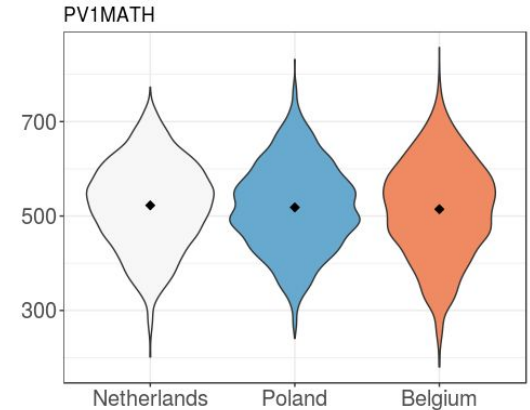Given a factor $C$ (with $k^*$ levels) and a numeric response $y^*$ analyze the differences among group means of $y$.

| | PV1MATH | CNT |
|---|---|---|
| 1 | 427.9561 | Belgium |
| 2 | 411.5984 | Poland |
| 3 | 471.1092 | Poland |
| 4 | 526.8032 | Netherlands |
| 5 | 721.4597 | Poland |
| 6 | 540.9020 | Poland |

# Problem 1

Given a factor $C$ (with $k$* levels) and
a numeric response $y$* analyze the differences
among group means of $y$.

| | PV1MATH | CNT | |
|---|---|---|---|
| 1 | 427.9561 | Belgium | |
| 2 | 411.5984 | Poland | |
| 3 | 471.1092 | Poland | |
| 4 | 526.8032 | Netherlands | |
| 5 | 721.4597 | Poland | |
| 6 | 540.9020 | Poland | |

PV1MATH



* k is greater than 2,
* y is normally distributed.

# Problem 1

Given a factor $C$ (with $k^*$ levels) and a numeric response $y^*$ analyze the differences among group means of $y$.

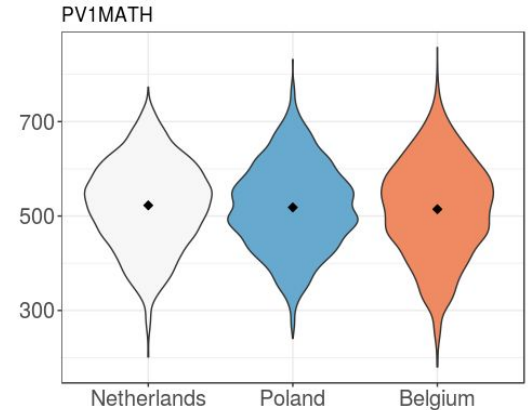# Solution

That's easy!

Let's run ANOVA and then post-hoc tests.

* k is greater than 2,
* y is normally distributed.

| | PV1MATH | CNT |
|---|---|---|
| 1 | 427.9561 | Belgium |
| 2 | 411.5984 | Poland |
| 3 | 471.1092 | Poland |
| 4 | 526.8032 | Netherlands |
| 5 | 721.4597 | Poland |
| 6 | 540.9020 | Poland |



PV1MATH

# Solution

Let's run ANOVA and then post-hoc tests.

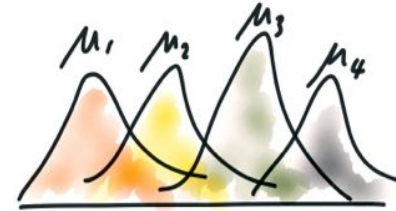# Solution

Let's run ANOVA and then post-hoc tests.

Let's try this out.

## Solution

Let's run ANOVA and then post-hoc tests.



Group means
with 95% confidence intervals

```
pisaAOV <- aov(PV1MATH ~ CNT, pisaNPB)
summary(pisaAOV)
#>               Df    Sum Sq Mean Sq F value  Pr(>F)
#> CNT            2     84272   42136   4.836 0.00796 **
#> Residuals  11113  96829278    8713
#> ---
#> Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1
```

# Solution

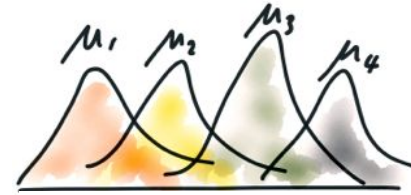Let's run ANOVA and then post-hoc tests.



95% family-wise confidence level



```
TukeyHSD(pisaAOV)
#> $CNT
#>                            diff       lwr        upr      p adj
#> Netherlands-Belgium   8.168042  1.827307 14.5087766  0.0071606
#> Poland-Belgium        3.793268 -1.962187  9.5487229  0.2700258
#> Poland-Netherlands   -4.374774 -9.166789  0.4172409  0.0819931
```

# Problem 2

Given a factor $C$ (with $k*$ levels) and
a numeric response $y*$ analyze the differences
among group means of $y$.

Group levels of $C$ into non-overlapping clusters.

95% family-wise confidence level

Differences in mean levels of Country

95% family-wise confidence level

Potentially $\binom{n}{3}$ inconsistencies

Differences in mean levels of Country

95% family-wise confidence level

Potentially $\binom{n}{3}$ inconsistencies

Cluttered visualizations

Differences in mean levels of Country

95% family-wise confidence level

Potentially $\binom{n}{3}$ inconsistencies

Cluttered visualizations

Fixed significance level

Differences in mean levels of Country

# Time for
factorMerger

# PISA 2012
Results in mathematics by country



## Group means
with 95% confidence intervals

(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood

GIC penalty = 12.5

1101991

1099551
1099376

## ANOVA table

|        | Df    | F        | p-value    |
|--------|-------|----------|------------|
| factor | 25    | 105939.6 | < 2.2e-16  |
| Res    | 92411 |          |            |

# Working with factorMerger

# Merge

1. Likelihood Ratio Tests
2. Delete or Merge Regressors

```
factorMerger::mergeFactors(response = myResponse,
                factor = myFactor,
                method = "LRT")
```

```
factorMerger::mergeFactors(response = myResponse,
                factor = myFactor,
                method = "hclust",
                successive = TRUE)
```

# Merge

1. Likelihood Ratio Tests
2. Delete or Merge Regressors

---

**Algorithm 1** Merging with $LRT$

---

    **function** MERGEFACTORS($response, factor, successive$)

2:      $pairsSet := generatePairs(response, factor, successive)$

       $M_0 :=$ full model

4:      **while** $levels(factor) > 1$ **do**

          $toBeMerged := \mathrm{argmax}_{pair \in pairsSet} l(updateModel(M_0, pair))$

6:         $M_0 := updateModel(M_0, toBeMerged)$

          $factor := mergeLevels(factor, pair)$

8:         $pairsSet := pairsSet \setminus pair$

      **end while**

10: **end function**

---

# Merge

1. Likelihood Ratio Tests
2. Delete or Merge Regressors

**Algorithm 2** Merging with agglomerative clustering

    **function** MERGEFACTORS(*response, factor, successive*)
2:        *pairsSet* := *generatePairs(response, factor, successive)*
        *dist* := set of distances
4:        **for all** *pair* $\in$ *pairsSet* **do**
           $h := \{\mu_{pair_1} = \mu_{pair_2}\}$                         $\triangleright$ hypothesis under which *pair* is merged
6:           $dist[pair] = LRT(M_h|M_0)$
        **end for**
8:        **if** successive **then**
           $hClust(dist,$ method = "single")
10:      **else**
           $hClust(dist,$ method = "complete")
12:      **end if**
    **end function**

More about the DMR algorithm: *https://arxiv.org/abs/1505.04008*

# PISA 2012
## Results in mathematics by country



(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood

-550763    -550532    -550302    -550071    -549841

# PISA 2012
Results in mathematics by country



A cluster

(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

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

Group means

Model's likelihood

-550763    -550532    -550302    -550071    -549841
loglikelihood

# PISA 2012
Results in mathematics by country



(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

LRT vs. full model

1e-70      1e-50      1e-30      1e-10

p-value

# PISA 2012
Results in mathematics by country



(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood

GIC penalty

GIC penalty = 12.5

Models:
constant, full
best

1101991

1099551
1099376

# PISA 2012

Results in mathematics by country



(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

-550763    -550532    -550302    -550071    -549841

logLikelihood

GIC penalty = 2

1101979

1099268

# PISA 2012

Results in mathematics by country



(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

logLikelihood

-550763   -550532   -550302   -550071   -549841

GIC penalty = 200

1104238

1102178

1100310

**PISA 2012**
Results in mathematics by country

**Group means**
with 95% confidence intervals

(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood: -550763  -550532  -550302  -550071  -549841

400  450  500  55

**PISA 2012**
Results in mathematics by country

**Groups frequencies**

(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood: -550763  -550532  -550302  -550071  -549841

0  5000  10000  15000

**PISA 2012**
Results in mathematics by country

**Boxplot**
Summary statistic: mean

(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood: -550763  -550532  -550302  -550071  -549841

200  400  600  800

**PISA 2012**
Results in mathematics by country

**Tukey HSD test**

(Estn): 528.48
(Nthr): 522.63
(Fnln): 520.72
(Plnd): 518.25
(Blgm): 514.46
(Grmn): 514.03
(Astr): 503.89
(Irln): 501.82
(Slvn): 500.32
(Dnmr): 498.82
(CzcR): 497.63
(UntK): 496.06
(Frnc): 495.94
(Prtg): 490.2
(Nrwy): 488.91
(Itly): 486.64
(Span): 482.16
(RssF): 482.13
(SlvR): 479.88
(Swdn): 479.85
(Hngr): 477.82
(Crot): 471.21
(Serb): 451.83
(Blgr): 439.89
(Mntn): 405.21

loglikelihood: -550763  -550532  -550071  -549841

a  b  c  d  e  f  g  h

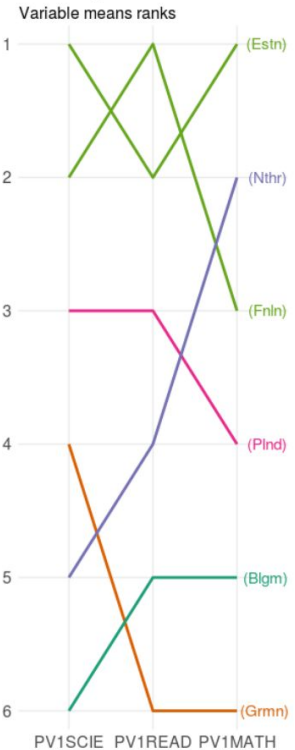# Other parametric models

1. multi dimensional Gaussian model,
2. binomial model,
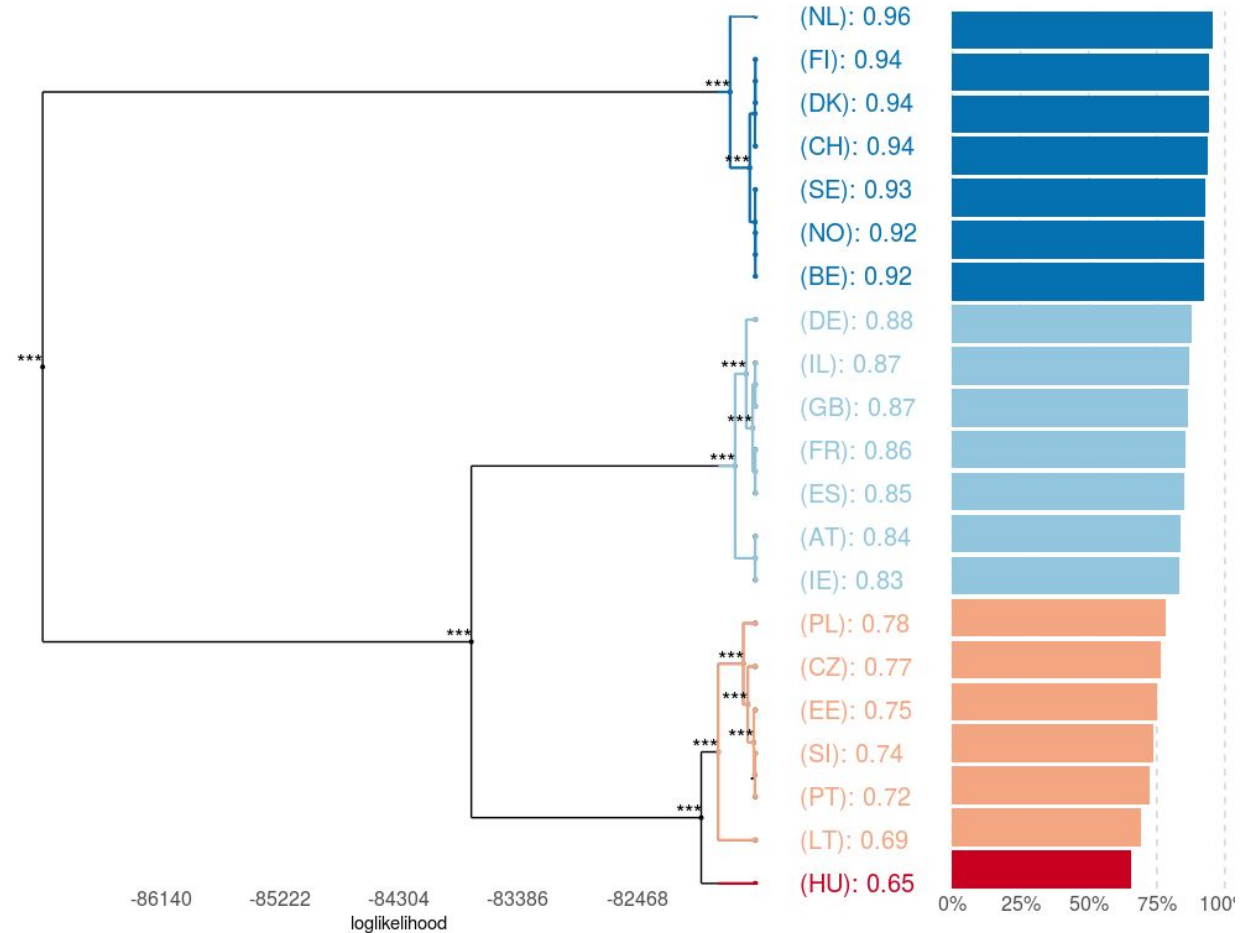3. survival model.



PISA 2012 - students' performance

(Estn): 23.36
(Fnln): 21.53

(Plnd): -2.32

(Nthr): -7.24
(Grmn): -9.07

(Blgm): -26.26

-395138    -395042    -394946    -394849    -394753
loglikelihood

Profile plot
Variable means ranks

(Estn)
(Nthr)
(Fnln)
(Plnd)
(Blgm)
(Grmn)

PV1SCIE  PV1READ  PV1MATH

## Other parametric models

1. multi dimensional Gaussian model,
2. binomial model,
3. survival model.



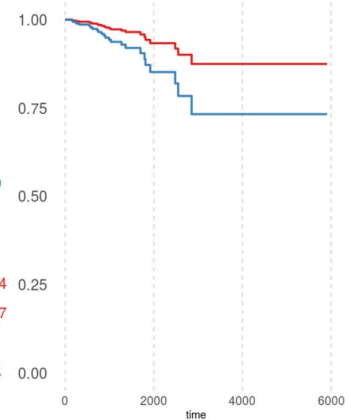European Social Survey 2014 - HOW HAPPY ARE YOU?    Success ratio

(NL): 0.96
(FI): 0.94
(DK): 0.94
(CH): 0.94
(SE): 0.93
(NO): 0.92
(BE): 0.92
(DE): 0.88
(IL): 0.87
(GB): 0.87
(FR): 0.86
(ES): 0.85
(AT): 0.84
(IE): 0.83
(PL): 0.78
(CZ): 0.77
(EE): 0.75
(SI): 0.74
(PT): 0.72
(LT): 0.69
(HU): 0.65

-86140   -85222   -84304   -83386   -82468
loglikelihood

0%   25%   50%   75%   100%

# Other parametric models

1. multi dimensional Gaussian model,
2. binomial model,
3. survival model.



Breast Cancer - treatment

(txtr): 2.35
(cytx): 1.72
(tmxf): 1.39
(Othr): 1.2
(armd): 0.84
(adrm): 0.67
(cycl): 0.47
(dxrb): 0.34

-131

loglikelihood

GIC penalty = 2

Survival plot
Adjusted survival curves for coxph model

time

ANOVA table

| | loglik | Chisq | Df | p-value |
|---|---|---|---|---|
| NULL | -131.2 | | | |
| factor | -128.4 | 5.4 | 7 | 0.606233 |

# Install and use the package

```r
install.packages("factorMerger")
```

CRAN

```r
if (!require(devtools)) install.packages("devtools")
devtools::install_github("geneticsMiNIng/factorMerger")
```

Github

```r
library(factorMerger)
fm <- mergeFactors(response = myResponse,
                   factor = myFactor,
                   family = "survival",
                   successive = TRUE,
                   method = "LRT")
```

Find more: *https://github.com/geneticsMiNIng/factorMerger*

# The aim of the factorMerger package

1. Create an algorithm which outputs an unequivocal data partition.
2. Improve visualizations.
3. Include other parametric models:
   a. multi dimensional Gaussian model,
   b. binomial model,
   c. survival model.

# geneticsMiNIng

# Any questions?

Agnieszka Sitko

ag.agnieszka.sitko@gmail.com

06-07-2017