

Mapping SAPA data

PMC lab
Northwestern University

Abstract

A guided tour in using the SAPA mapping functions.

To graphically display the SAPA data can be done using a number of external packages as well as a set of functions developed at the PMC lab. This is a short guide to using these sapa functions. Comments are welcome. Suggestions for improvement of the functions are particularly welcome.

Install other packages and then make them active

The most relevant package is the *zipcode* and its associated *zip code* database. In addition, the *mgcv* package is used for finding those points within a map.

installing the relevant packages

```
install.packages("zipcode")
install.packages("mgcv")
```

Activate them

Then, once they are installed, they need to be activated every time you start R.

```
library(zipcode)
library(mgcv)
```

Get the SAPA mapping functions

The current working version of the SAPA mapping functions may be sourced from the personality-project.org site. Eventually these should be wrapped into the package that Jason and David are developing at ICAR-project.org

```
source("http://personality-project.org/r/src/sapa.r")
```

contact: William Revelle revelle@northwestern.edu
Draft version of October 7, 2013
Please do not cite without permission

Using the SAPA mapping function

Getting the zip codes attached to the summary statistics

The `sapa.zip` takes a data file (e.g., `IRT.scores`) and combines it with the zip codes from the zip code database. The output is a data.frame of the means and ns of the original variables by zip code.

This just needs to be done once for a particular data set. Maps for particular dependent variables are drawn using this data base.

```

s.zip <- sapa.zip(IRT.scores)
dim(s.zip)
colnames(s.zip)

3208 100

[1] "zip"
[4] "age"
[7] "BMI"
[10] "country"
[13] "education"
[16] "jobstatus"
[19] "p1occ"
[22] "p2edu"
[25] "Boldness"
[28] "Honesty"
[31] "Intellect"
[34] "Politeness"
[37] "Extraversion"
[40] "EmotionalStability"
[43] "RightWingAuthoritarianism"
[46] "MXiq"
[49] "gender.n"
[52] "marstatus.n"
[55] "ParentalEdu.n"
[58] "state.n"
[61] "education.n"
[64] "jobstatus.n"
[67] "p1occ.n"
[70] "p2edu.n"
[73] "Boldness.n"
[76] "Honesty.n"
[79] "Intellect.n"
[82] "Politeness.n"
[85] "Extraversion.n"
[88] "EmotionalStability.n"
[91] "RightWingAuthoritarianism.n"
[94] "MXiq.n"
[97] "city"
[100] "longitude"

[1] "gender"
[4] "marstatus"
[7] "ParentalEdu"
[10] "state.x"
[13] "discipline"
[16] "jobfield"
[19] "p1edu"
[22] "Assertiveness"
[25] "Compassion"
[28] "Humility"
[31] "Openness"
[34] "Agreeableness"
[37] "Integrity"
[40] "GunControl"
[43] "g"
[46] "R3Diq"
[49] "relstatus.n"
[52] "height.n"
[55] "weight.n"
[58] "ethnic.n"
[61] "discipline.n"
[64] "jobfield.n"
[67] "p1edu.n"
[70] "Assertiveness.n"
[73] "Compassion.n"
[76] "Humility.n"
[79] "Openness.n"
[82] "Agreeableness.n"
[85] "Integrity.n"
[88] "GunControl.n"
[91] "g.n"
[94] "R3Diq.n"
[97] "state.y"
[100] "relstatus"
[4] "height"
[7] "weight"
[10] "ethnic"
[13] "major"
[16] "occupation"
[19] "p2occ"
[22] "Balance"
[25] "Enthusiasm"
[28] "Industriousness"
[31] "Orderliness"
[34] "Conscientiousness"
[37] "OpennessNewExp"
[40] "Environmental"
[43] "ANiq"
[46] "VEdiq"
[49] "age.n"
[52] "BMI.n"
[55] "country.n"
[58] "zip.n"
[61] "major.n"
[64] "occupation.n"
[67] "p2occ.n"
[70] "Balance.n"
[73] "Enthusiasm.n"
[76] "Industriousness.n"
[79] "Orderliness.n"
[82] "Conscientiousness.n"
[85] "OpennessNewExp.n"
[88] "Environmental.n"
[91] "ANiq.n"
[94] "VEdiq.n"
[97] "latitude"

```

Create a longitude by latitude grid of the data

This next step creates a grid of the means for each zip code for a particular Dependent Variable (e.g., 'g'). Data that are in the same grid are averaged. Note that the larger

the grid size, the more precise the results, but the time taken for this and subsequent calculations increase by the square of the grid dimensions. There is an 'average' option which subtracts the overall weighted average from every cell. This functionally centers the data around that average. The alternative is to have all empty cells have the average.

```
iq.grid <- sapa.grid(data=s.zip,DV='g',grid=300,average=TRUE)
names(iq.grid)
dim(iq.grid$av)

[1] "average" "total"   "count"

[1] 300 300
```

Smooth the data

Although not necessary for graphics, smoothing the data makes an overall more pleasing map. The smooth function takes a normal “kernel” and applies it to every cell. That is, nearby cells are weighted as a function of the density of the normal curve. The default value of 9 means that the kernel is a 9 x 9 grid. Each grid step is taken to be half a sd of the normal curve. The miss parameter specifies how much to weight cells with no observations. By increasing the weight, we are increasing the tendency to plot average values for missing cells.

The smoothing function operates on all cells of the grid.

It is useful to compare smoothed data with the raw plots.

```
iq.smooth <- sapa.smooth(iq.grid,size=11,miss=.05)
dim(iq.smooth)

300 300
```

Find the points inside the map.

The procedure for smoothing that I am using will create points outside the map region. This is rectified by using a function from the *mvcv* package, *inSide*. In addition, the *sapa.inside* function converts the grid based data to a list of x, y, and z coordinates that is used by the *sapa.plot* function.

For some reason the *inSide* function gives a warning message which I think can be ignored.

In the example below, I do this operation for both the smoothed data and the raw data. This will allow for a comparison of the two approaches. The *sapa.inside* returns a list of x, y, and z for plotting.

```
iq.smoothed.inside <- sapa.inside(iq.smooth)

iq.inside <- sapa.inside(iq.grid$average)

names(iq.smoothed.inside)

"x" "y" "z"
```

Now plot the results

We can either plot the raw data points (unsmoothed) or the smoothed data points. Even better, we can now draw image plots of the smoothed or unsmoothed data. These image plots make for much prettier pictures.

This one was done using a 400 x 400 grid and 51 colors.

```
sapa.plot(iq.inside)
sapa.image(iq.smooth.inside,main="iq by state",n=51,database="state")
```

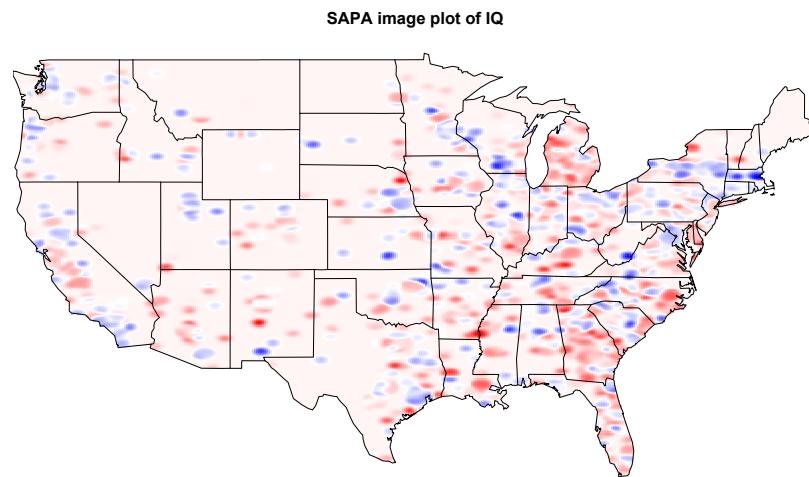


Figure 1. IQ by state, smoothed

And One function to Rule them All

It is obviously tedious to enter the four lines necessary to make pretty pictures, so, lets try the **sapa.combined** function. You must first run the **sapa.zip** function once to prepare the data base. Then a call to **sapa.combined** will do the rest. The function returns the object for remapping with different parameters.

```
my.zips <- sapa.zip(IRT.scores)
g <- sapa.combined(my.zips,'g',gridsize=400,main='ICAR IQ scores from the SAPA project')
```

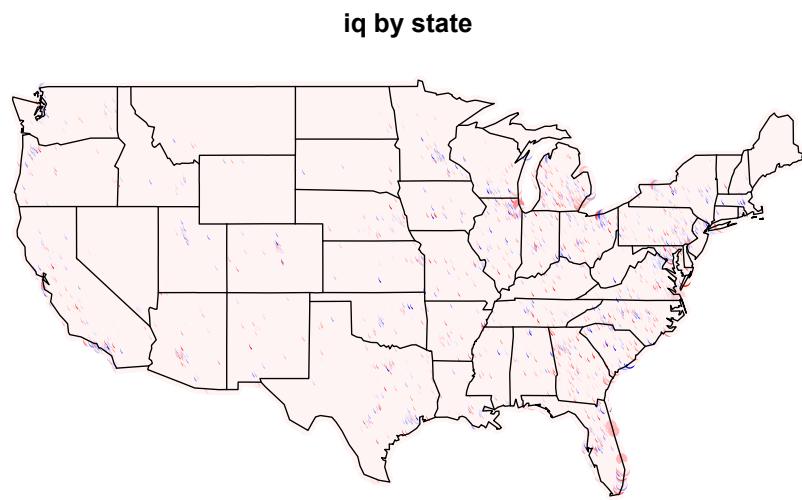


Figure 2. IQ by state, raw

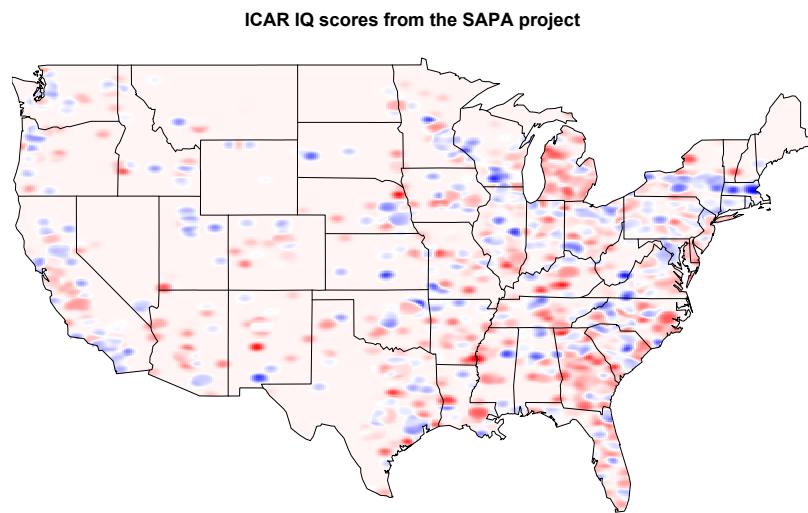


Figure 3. IQ from ICAR for the US. Smoothed size = 11 and with a grid size of 400. By default, the number of colors is the grid size (400).