mixtureReg: A Quick Start

In this tutorial, we are going to show how to use **mixtureReg** to model data with two possible regimes.

Data

The data used for demostration purpose here is the CO2 data set from the mixtools package.

```
library(mixtools)
## mixtools package, version 1.0.4, Released 2016-01-11
## This package is based upon work supported by the National Science Foundation under Grant No. SES-051
data("CO2data")
head(CO2data)
       GNP CO2 country
## 1 19.02 14.7
                    CAN
## 2 3.67 3.9
                    MEX
## 3 28.20 20.8
                    USA
## 4 40.94 9.0
                    JAP
## 5 10.61 8.3
                    KOR
```

A simple example

AUS

6 20.09 16.0

The motivation of mixture of regressions is that there can be two different regimes in the data so we want to fit two lines through the data.

We can easily achieve this by putting two regression formula into a list and feed it into the **mixtureReg** function.

In this case, the message shows that the model converges in 32 iterations.

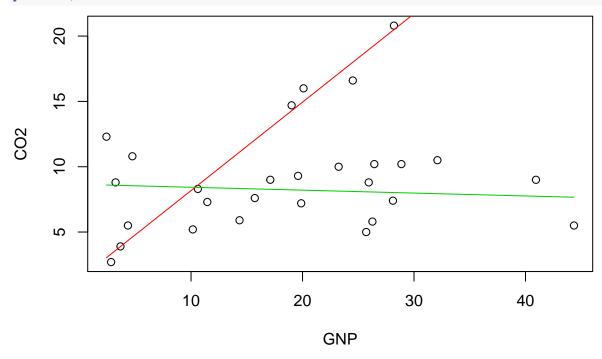
```
library(mixtureReg)
mx1 <- mixtureReg(</pre>
 regData = CO2data,
  formulaList = list(formula(CO2 ~ GNP),
                      formula(CO2 ~ GNP))
)
## Loading required package: dplyr
##
## Attaching package: 'dplyr'
## The following objects are masked from 'package:stats':
##
##
       filter, lag
## The following objects are masked from 'package:base':
##
       intersect, setdiff, setequal, union
##
```

```
## diff = 5.214417e-09
## iter = 31
## restart = 0
## log-likelihood = -66.98373
```

The fit

We provide a plot method (S3 method) to visualize the predictions from the model. The circles below are the original data points and the red lines are predictions from our model.

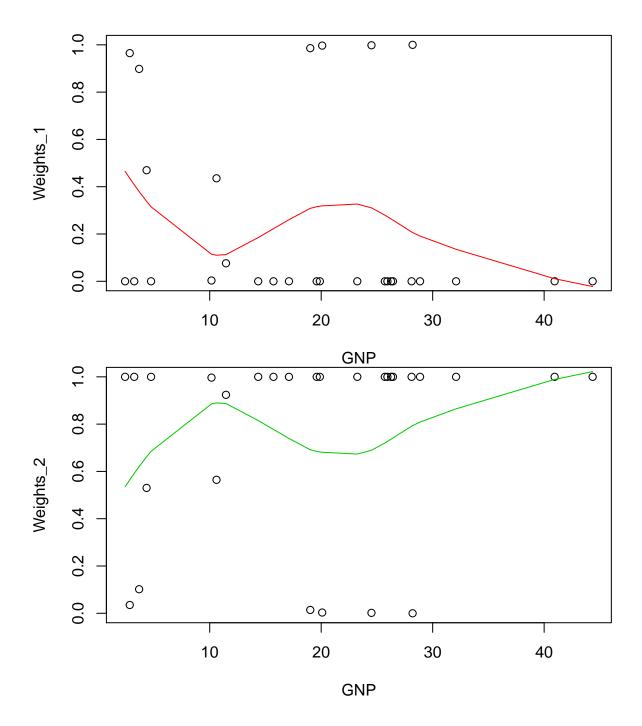
plot(mx1, which = 1)



The weights

Other than predictions, The mixture of regressions also produces weight estimates for each data point which indicate the posterior probabilities of membership to the regression lines. We provide another plot method to visualize these weights.

plot(mx1, which = 2)



The iterations

(Mainly for debugging purposes) We also provide a monitor component for modelers to learn more about what are happening in iterations.

head(mx1\$monitor)

```
diff iter restart
                               logLik
                                           newLL
                                                   sigma1
                                                            sigma2
                                                                       ratio
## 1 1.000000
                  0
                          0 -77.79608
                                              NA 2.980934 2.735892 1.089565
## 2 0.4036337
                  1
                          0 -77.39245 -77.39245 3.092703 2.557185 1.209417
## 3 0.8106541
                  2
                          0 -76.58179 -76.58179 3.184343 2.295055 1.387480
## 4 0.9972444
                  3
                          0 -75.58455 -75.58455 3.178627 2.032839 1.563640
```

```
## 5 0.7996175
                           0 -74.78493 -74.78493 3.066775 1.855304 1.652977
## 6 0.4896748
                  5
                           0 -74.29526 -74.29526 2.894154 1.784004 1.622280
       lambda1
##
                 lambda2 error message
## 1 0.4822502 0.5177498
                                     NA
## 2 0.4792906 0.5207094
                                     NA
## 3 0.4705942 0.5294058
                                     NA
## 4 0.4518643 0.5481357
                                     NA
## 5 0.4250229 0.5749771
                                     NA
## 6 0.3958220 0.6041780
                                     NA
```

Flexible modeling

A nice feature of this package is that we can flexibly specify the formula as we would in lm.

For example, we can restrict one regression line to be horizontal with no slope coefficient.

This example also helps to demonstrate the usage of "yName" and "xName" arguments in the plot method. When not specified, the plot method will search the variables in the first formula, which will not work in this case.

```
mx2 <- mixtureReg(</pre>
 regData = CO2data,
  formulaList = list(formula(CO2 ~ 1),
                      formula(CO2 ~ GNP))
## diff =
           3.496524e-09
## iter =
           25
## restart =
## log-likelihood = -67.10555
plot(mx2, yName = "CO2", xName = "GNP", which = 1)
     20
                                                  0
             0
                 0
                                                               0
     10
                                                     0
                                                          0
                                                0
                                          0
                                      0
                                                                              0
                                                     0
              0
                                   0
                                                        0
                            0
                                          0
                                                    0
                                 0
                                                                                    0
                          0
      S
                         10
                                          20
                                                           30
                                                                            40
```

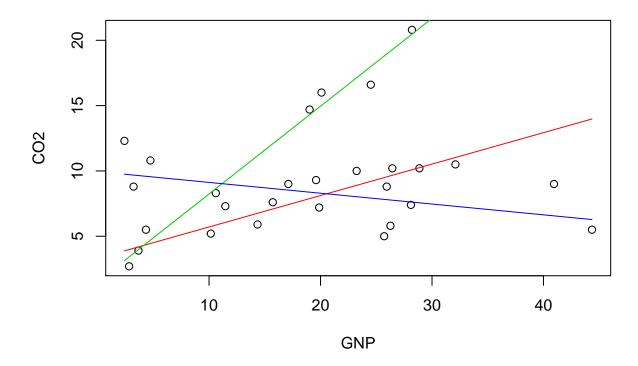
We can also specify 2nd order polynomial lines.

GNP

```
mx3 <- mixtureReg(</pre>
  regData = CO2data,
  formulaList = list(formula(CO2 ~ GNP + I(GNP^2)),
                     formula(CO2 ~ GNP + I(GNP^2)))
## diff = 0.001268958
## iter = 45
## restart = 15
## log-likelihood = -65.42241
plot(mx3, which = 1)
     15
             0
                                                             0
     10
                                                    0
                                                      0
                                               0
                                                                            0
              0
                          0
                                                       0
                                         0
                                0
                                                                                 0
     2
               6
             0
                         10
                                         20
                                                         30
                                                                          40
```

We can also fit three (or more) regressions at the same time.

GNP



Comparison with mixtools

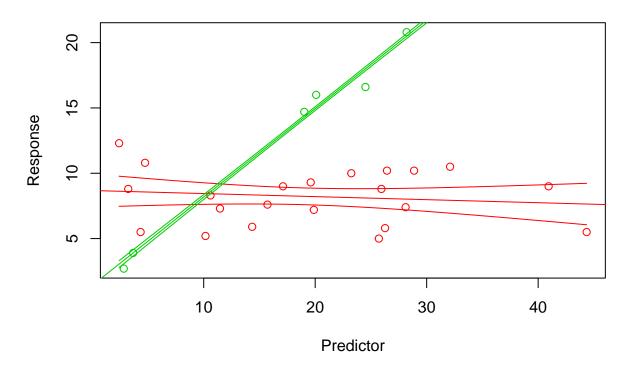
The main shortcoming of **mixtools** is that it doesn't provide easy to use options to restrict model coefficients like we do in model 2.

The following is an example from Tatiana Benaglia, Didier Chauveau, David R. Hunter, Derek Young (2009). This example produces similar results with our model 1.

```
compare1 <- mixtools::regmixEM(
   CO2data$CO2, CO2data$GNP,
   lambda = c(1/4, 3/4),
   beta = matrix(c(2, 0, 0, 1), 2, 2),
   sigma = c(1,1)
   )</pre>
```

```
## number of iterations= 18
plot(compare1, whichplots = 2)
```

Most Probable Component Membership



References

de Veaux RD (1989). "Mixtures of Linear Regressions." Computational Statistics and Data Analysis, 8, 227-245.

Tatiana Benaglia, Didier Chauveau, David R. Hunter, Derek Young (2009). mixtools: An R Package for Analyzing Finite Mixture Models. Journal of Statistical Software, 32(6), 1-29. URL http://www.jstatsoft.org/v32/i06/.