

Parameter Estimation of Compartment Models in **SoilR** Using Classical and Bayesian Optimization

Markus Müller* and Carlos A. Sierra†
Max Planck Institute for
Biogeochemistry

December 19, 2019

1 Introduction

The objective of this document is to provide examples on how to use **SoilR** in combination with package **FME** to infer parameter values of soil organic matter decomposition models using observed data. Parameter estimation for dynamical systems is an advanced topic of inverse modeling and as such is far beyond the scope of this vignette. We will point to some principal questions and possible problems as they arise, but this treatment will be far from comprehensive. This document also does not replace the documentation of package **FME** (Soetaert & Petzoldt, 2010), which we strongly recommend to consult. Instead, we show first how a small wrapper function makes a **SoilR** model available for the functions in **FME**, and second, how to choose the right parameterizations of **SoilR** models to meet the requirements of the **FME** algorithms. We present two examples. One is the parameterization of a two-pool model applied to a soil incubation experiment. The other example uses observed radiocarbon data from CO₂ measurements conducted at Harvard Forest, USA.

2 First Example: A soil incubation experiment

Measurements of evolved CO₂ from incubation experiments can provide useful data for parameterizing soil organic matter decomposition models and identify functionally distinct pools (Schadel et al., 2013). We present here data from an incubation experiment in which we measured the evolved CO₂ from a forest soil. The dataset `incubation_experiment`, is already included in **SoilR** and contains data from an incubation experiment with a boreal forest soil. After loading **SoilR** into our R session we can explore its parts, extract the flux data from the boreal site into a separate object and plot it.

*mamueller@bgc-jena.mpg.de

†csierra@bgc-jena.mpg.de

```

library(SoilR)
summary(incubation_experiment)

##              Length Class      Mode
## eC02          3      data.frame list
## c_concentrations 3      -none-    numeric
## soil_mass       1      -none-    numeric

BorealC02=incubation_experiment$eC02
plot(
  BorealC02[,1:2]
  , xlab="Days"
  ,ylab="CO2 flux in (mg C /g_soil/day)"
  ,ylim=c(0,50)
)
arrows(BorealC02[,1],BorealC02[,2]-BorealC02[,3],BorealC02[,1],
       BorealC02[,2]+BorealC02[,3],code=3,angle=90,length=0.1)

```

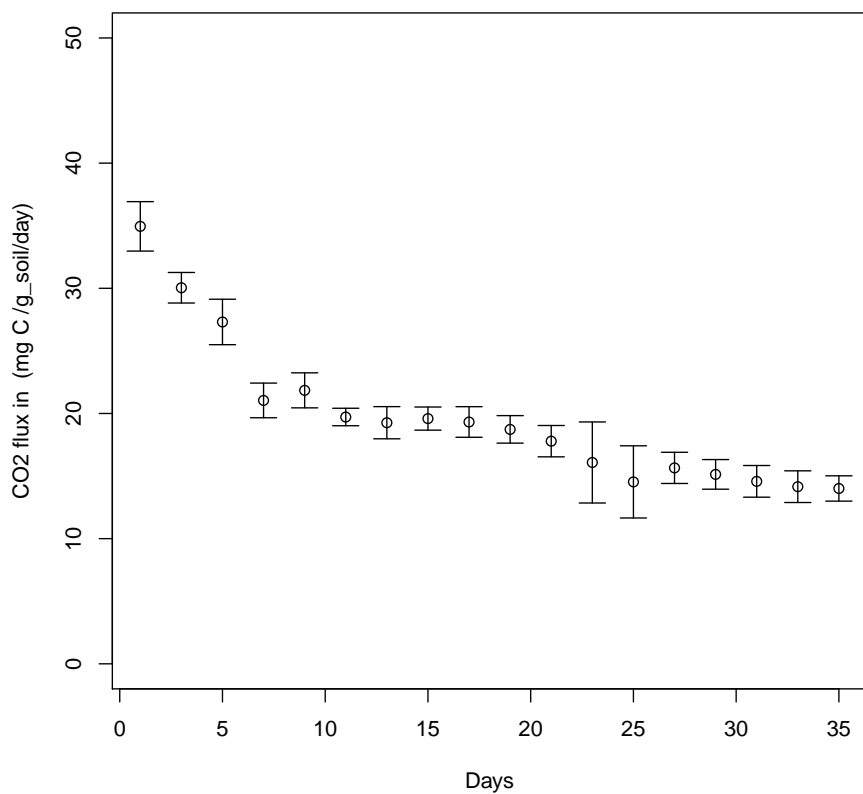


Figure 1: Respired CO₂ from an incubation experiment with a boreal forest soil.

The models in `SoilR` do not use units, but rather assume that the input data is consistent in this respect. If we want to supply the initial mass in *g* and the time in days, then we have to provide the flux in *g* per day, whereas the original data is given in mg carbon per g of soil per day. We therefore rescale the columns for the fluxes and the error.

```
Ctotal<- mean(incubation_experiment$c_concentrations) *
           incubation_experiment$soil_mass
BorealCO2<-data.frame(
  BorealCO2[,1]
  ,BorealCO2[,2:3]*1e-06*Ctotal)

# rename the columns
names(BorealCO2)<-c("time", "eCO2", "eCO2sd")
# create a function for later plots
plot_with_data<-function(x){
  plot(
    x=x
    ,xlab="Days"
    ,ylab="Evolved CO2 (gC/day)"
    ,ylim=c(0 ,9e-04)
    ,type='l'
  )
  points(
    BorealCO2
  )
  arrows(
    BorealCO2[,1]
    ,BorealCO2[,2]-BorealCO2[,3]
    ,BorealCO2[,1]
    ,BorealCO2[,2]+BorealCO2[,3]
    ,code=3
    ,angle=90
    ,length=0.1
  )
}
```

Before we embroil ourselves in technicalities we should think about the general possibility to identify the parameters from a single time line of the combined release of a yet unknown number of pools with yet unknown connections and outlets. We face several challenges, in particular we have to:

1. Find a class of models that is large enough to contain a model that is capable of reproducing the observed data. If the class is too small, we will “underfit” our data and the model will be “biased”.
2. Find a number of parameters that is small enough to be actually determined unambiguously from the data. If we have parameters whose effects enhance or cancel the effects of other parameters, the parameterized model will be “overfitted” and make (possibly extremely) misleading predictions

for data not included in the training set. (In our case the predicted timeline of an overfitted model could meet all the measurements very well but be completely unreasonable in between or beyond the measurement times.)

3. Make sure that the parameters can at least be tested independently by the optimization procedure. (We refer here not to the desired independence of the parameters w.r.t (with respect to) the impact on the result as mentioned under 2., but to the shape of the set of valid parameters. For n parameters the algorithms that we will use accept n ranges and assume to be able to choose *any* combination of parameters out of this n -dimensional rectangular set to form a *valid* model. For example, an optimization algorithm that changes parameters independently should not accidentally create a model that breaks mass conservation or produces negative fluxes. `SoilR` allows many ways to specify models, including functions whose arguments include matrices whose parameters are *not* independent. Since it always checks that the specified model is a valid compartmental system, it would actually stop the optimization procedure from trying parameters that can not possibly constitute a soil model. We will choose parameterizations that avoid this situation.)

We will start with a simple model that fulfills conditions 2. and 3. and extend it to improve the degree of accuracy until we can no longer guarantee 2. Note that the search for accurate and identifiable soil models is an open research question. A recent study Sierra et al. (2015) suggests that for the possibilities to uniquely identify (structural) parameters from incubation data are severely limited. We can usually only hope to identify 2 or maybe 3 parameters. Our first attempt to model the data is therefore a very modest one pool model with a constant decomposition rate k .

2.1 One pool, one parameter model

The first FME function that we want to use is `modFit`. A look at its documentation (type `? modFit`) suggests that we at least have to provide a function `f` to be minimized, a vector of start values for the parameters, and in our case, a lower limit for the parameter k since we know that a negative value would create an invalid (not mass conserving) model. The function `f` will be a “costfunction” that will actually need to evaluate the model as a function of the parameter(s) at the times where we have data and compare the output to the measured data. It will return a vector of residuals (one entry for every measurement time) possibly weighted by the accuracy of our data. To help us write such a function, FME provides the function `modCost` to automatically create the cost function that compares the result of our model run with our data and computes the residuals weighted by the measurement errors if available. What remains to be done for us is to create a function of the parameter(s) that produces the release flux for the times in our dataset.

```
library(FME)

## Loading required package: rootSolve
## Loading required package: coda
```

```
eC02P1=function(pars){
  At=ConstLinDecompOp(
    out_flux_rates=ConstantOutFluxRateList_by_PoolIndex(list("1"=pars[[1]]))
    ,numberOfPools =1
  )
  mod=GeneralModel(
    t=BorealC02[, 'time']
    ,A=At
    ,ivList=c(Ctotal)
    ,inputFluxes=0
    #,pass=TRUE
  )
  Rt<-getReleaseFlux(mod)
  return(
    data.frame(time=BorealC02[, 'time'], eC02=rowSums(Rt)))
}
```

Notice that our function, `eC02func`, requires a (for this first example one dimensional) vector of parameters `pars` with the values of the first decomposition rate in position 1. Our function returns a `data.frame` with two columns, time in days and the sum of the release flux for the two pools. The next step is to create a cost function according to `FME` requirements. This cost function takes as arguments a function with the model, the set of observations, and a measure of the error in the observations. The function calculates sums of squared residuals from the model output and the observed data, weighted by the standard deviation of the measurement. For convenience we use the `FME` function `modCost` which returns an object of a class with the same name.

```
eC02P1cost=function(pars){
  return(
    modCost(
      model=eC02P1(pars)
      ,obs=BorealC02
      ,err="eC02sd"
    )
  )
}
```

Having defined the cost function we only need initial values (`inivars`) and (optional) lower and upper bounds (`upP1`, `loP1`) for the parameters to finally use the function `modFit`. It implements the optimization as an iterative process. First, the cost function will be evaluated on the initial parameter values. Then, the algorithm will try to guess new parameters, evaluate the cost function again and repeat this process until the cost is small enough or the number of permitted iterations exceeded.

```
inipars=c( 1/2000 )
upP1=c(1)
```

```
loP1=c(0)

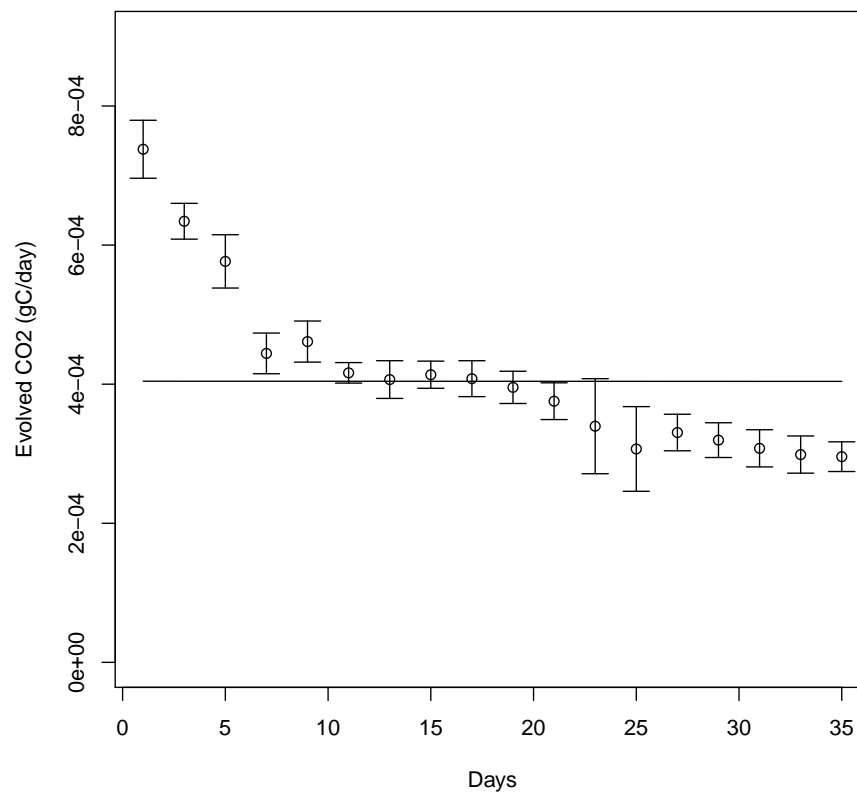
eC02P1fit=modFit(
  f=eC02P1cost
  ,p=inipars
  ,upper=upP1
  ,lower=loP1
)
```

To see the best set of parameter values found by the function we can type:

```
eC02P1fit$par
## [1] 1.91489e-05
```

This set of parameters can be used now to run the model again and plot the obtained model against the observations.

```
plot_with_data(eC02P1(eC02P1fit$par))
```



It is obvious that this model is too simple. The measured release flux is so small that - if all the carbon is decomposing - the decomposition rate must be also very small, indeed so small as to render the carbon almost stable. The expected exponential shape suggested by the data is not visible at all. We therefore probably can improve the fit significantly if we allow a second parameter determining how much of the total carbon is accessible for decomposition.

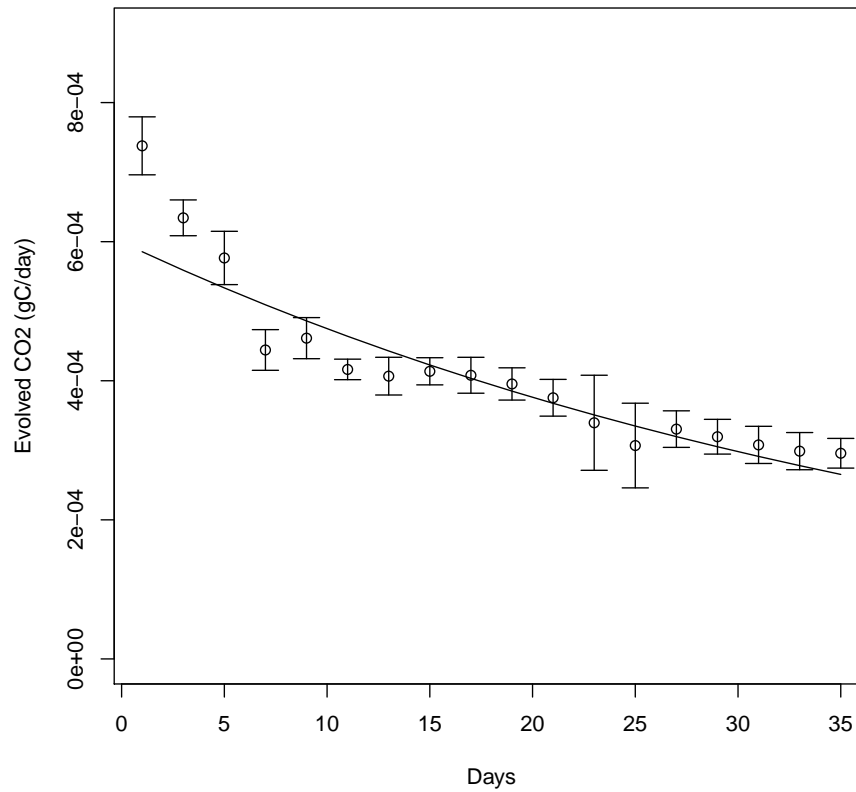
2.2 One pool, with two parameters

```
eC02P1a=function(pars){
  At=ConstLinDecompOp(
    out_flux_rates=ConstantOutFluxRateList_by_PoolIndex(list("1"=pars[[1]]))
    ,numberOfPools =1
  )
  mod=GeneralModel(
    t=BorealC02[, 'time']
    ,A=At
    ,ivList=c(Ctotal*pars[2])
    ,inputFluxes=0
  )
  Rt<-getReleaseFlux(mod)
  return(data.frame(time=BorealC02[, 'time'], eC02=rowSums(Rt)))
}

eC02P1acost=function(pars){
  return(
    modCost(
      model=eC02P1a(pars)
      ,obs=BorealC02
      ,err="eC02sd"
    )
  )
}

upP1a=c(1,1)
loP1a=c(0,0)

initPars=c(0.02, 0.01)
eC02P1afit=modFit(
  f=eC02P1acost
  ,p=initPars
  ,upper=upP1a
  ,lower=loP1a
)
plot_with_data(eC02P1a(eC02P1afit$par))
```



This time we also look at some estimated statistics.

```
eC02P1afit$par
## [1] 0.023284171 0.001191259
summary(eC02P1afit)
##
## Parameters:
##      Estimate Std. Error t value Pr(>|t|)
## [1,] 0.0232842  0.0027984   8.321 3.32e-07 ***
## [2,] 0.0011913  0.0001043  11.426 4.17e-09 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.683 on 16 degrees of freedom
##
## Parameter correlation:
##      [,1] [,2]
## [1,] 1.000 -0.959
## [2,] -0.959 1.000
```


According to these results we can say: Assuming a one pool model with constant decomposition rate with only a fraction of the total carbon accessible for decomposition, the best fitting model has a decomposition rate of 2.3% per day where only 0.12% of the total carbon are accessible for decomposition.

We can also check the estimated correlation between our two parameters. The high (anti)correlation indicates that it is difficult to identify the best model confidently since increasing the decomposition rate and decreasing the accessible fraction of carbon have similar effects on the model output (the release flux) .

The results of `modFit` can be used for Bayesian parameter estimation with `FME` 's Markov Chain Monte Carlo function `modMCMC`. In our example, we used the default values to avoid explanations about `FME` , but there are many possibilities to influence the result, which are briefly mentioned in the help (type ? `modMCMC`) and explained in more detail in Soetaert & Petzoldt (2010). To avoid long run-times, we only use 3000 iterations in this example, but this number can be much larger to guarantee convergence of the chains. The results of the MCMC procedure can be obtained with the function `summary()`. The output gives the mean, standard deviation, min and max, and 25% quantiles for all parameter values. It also produces summary statistics for the variance of the observed variable.

```
eC02P1amcmc=modMCMC(
  f=eC02P1acost
  ,p=eC02P1afit$par
  ,niter=3000
  ,upper=upP1a
  ,lower=loP1a
)

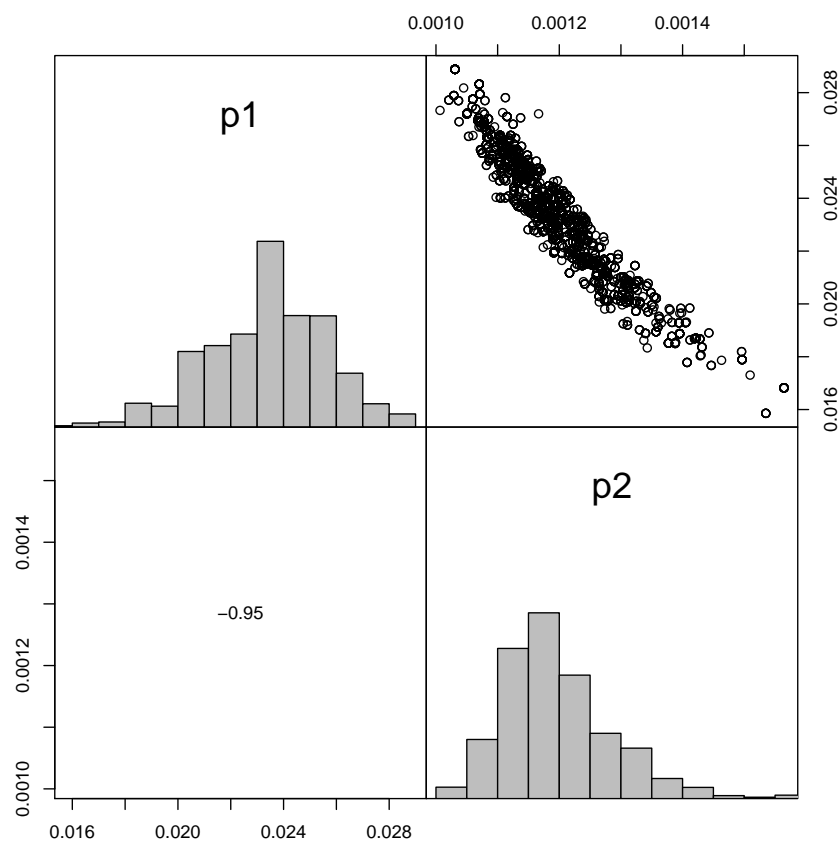
## number of accepted runs: 551 out of 3000 (18.36667%)

summary(eC02P1amcmc)

##              p1              p2
## mean 0.02338923 1.196762e-03
## sd   0.00228184 8.908731e-05
## min  0.01585660 1.006553e-03
## max  0.02888280 1.564456e-03
## q025 0.02186112 1.133328e-03
## q050 0.02354218 1.187631e-03
## q075 0.02500358 1.241460e-03
```

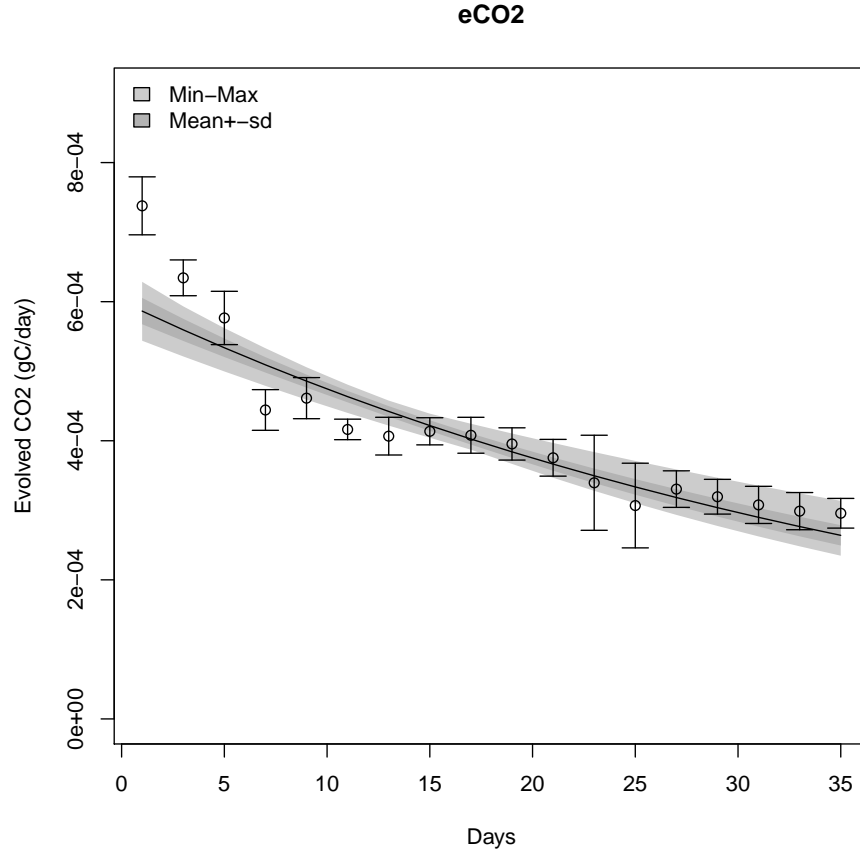
A plot with the posterior distribution of the obtained parameter values can be obtained with function `pairs`.

```
pairs(eC02P1amcmc)
```



The strong correlation between the parameters is again apparent. We can also visualize the results of the Monte Carlos sampling by estimating the effect on the results (at different times) with function `sensRange` and plotting the results as a function of time.

```
predRange=sensRange(func=eC02P1a, parInput=eC02P1amcmc$par)
plot_with_data(
  summary(predRange)
)
```



Note that the results depend not only on our equations (as they would if we had sampled the complete parameter space with infinitely close parameter pairs), but also on the the other arguments of `modMCMC`, which for instance determine how many parameter combinations will be tested or how much they differ, or in general how well the parameter space has been sampled.

2.3 Two pool model with three parameters

As a last example we try to further improve the fit of the model by allowing a second pool with its decomposition rate as second parameter. The role of the first parameter as constant decomposition rate for the first pool does not change. The second parameter of the last model becomes the third here and decides now how the total carbon is distributed between the two pools.

```
eCO2P2=function(pars){
  At=ConstLinDecompOp(
    out_flux_rates=ConstantOutFluxRateList_by_PoolIndex(list("1"=pars[[1]], "2"=pars[[2]]
    ,numberOfPools =2
  )
  gamma=pars[[3]]
  mod=GeneralModel(
```

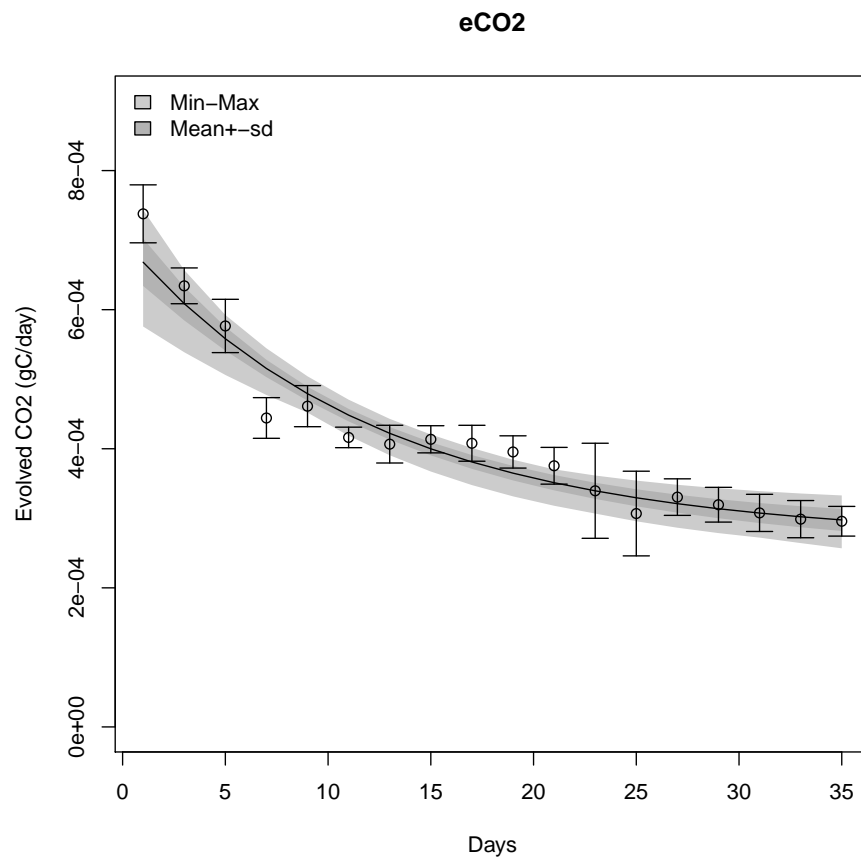
```

    t=BorealC02[, 'time']
    ,A=At
    ,ivList=Ctotal*c(gamma,1-gamma)
    ,inputFluxes=c(0,0)
    ,pass=TRUE
  )
  Rt<-getReleaseFlux(mod)
  return(data.frame(time=BorealC02[, 'time'], eC02=rowSums(Rt)))
}
#r=c(1:4, 34:35)
eC02P2cost=function(pars){
  return(
    modCost(
      model=eC02P2(pars)
      ,obs=BorealC02
      ,err="eC02sd"
    )
  )
}
initPars=c(.02, .0001, 0.999)
upP2=c(.1, .1, 1)
loP2=c(0, 0, 0)
eC02P2fit=modFit(
  f=eC02P2cost
  ,p=initPars
  ,method="Marq"
  ,upper=upP2
  ,lower=loP2
)
eC02P2mcmc=modMCMC(
  f=eC02P2cost
  ,p=eC02P2fit$par
  ,niter=3000
  ,upper=upP2
  ,lower=loP2
)

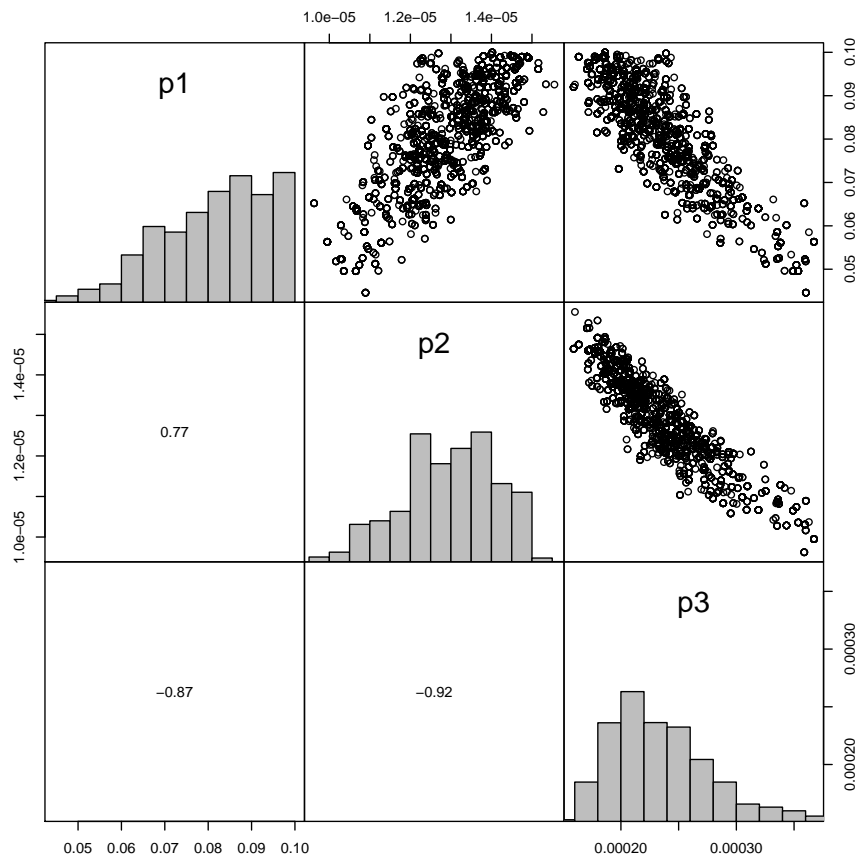
## number of accepted runs: 517 out of 3000 (17.23333%)

plot_with_data(
  summary(
    sensRange(
      func=eC02P2
      ,parInput=eC02P2mcmc$par
    )
  )
)

```



```
pairs(eCO2P2mcmc)
```



```
summary(ecO2P2fit)

##
## Parameters:
##      Estimate Std. Error t value Pr(>|t|)
## [1,] 9.998e-02 1.664e-02  6.010 2.39e-05 ***
## [2,] 1.402e-05 8.930e-07 15.699 1.02e-10 ***
## [3,] 1.862e-04 3.191e-05  5.837 3.27e-05 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 1.068 on 15 degrees of freedom
##
## Parameter correlation:
##      [,1] [,2] [,3]
## [1,] 1.0000 0.8478 -0.9124
## [2,] 0.8478 1.0000 -0.9384
## [3,] -0.9124 -0.9384 1.0000
```

We see again a high correlation between the parameters, which indicates poor

identifiability. Numerical experiments (e.g. variation of the startvalues) also show that the algorithms occasionally get lost in local minima, so that including more pools or fluxes does not seem to be a very good idea. Note however that although we probably cannot hope to fit models that are structurally more complex there are infinitely more ways to choose parameters. We could for instance define a time dependent decomposition rate whose shape depends on more than the one parameter we chose for a constant rate.

3 Second Example: Radiocarbon in respired CO₂

`SoilR` can also calculate the amount of radiocarbon in soils or in respired CO₂. Here, we take as an example a series of observations of radiocarbon in respired CO₂ conducted at Harvard Forest, USA. The dataset is also included in `SoilR`, and is visualized in Figure 2.

```
plot(  
  D14C~Year  
  ,data=HarvardForest14C02  
  ,ylab=expression(paste(Delta^14,"C",' per mille'))  
)
```

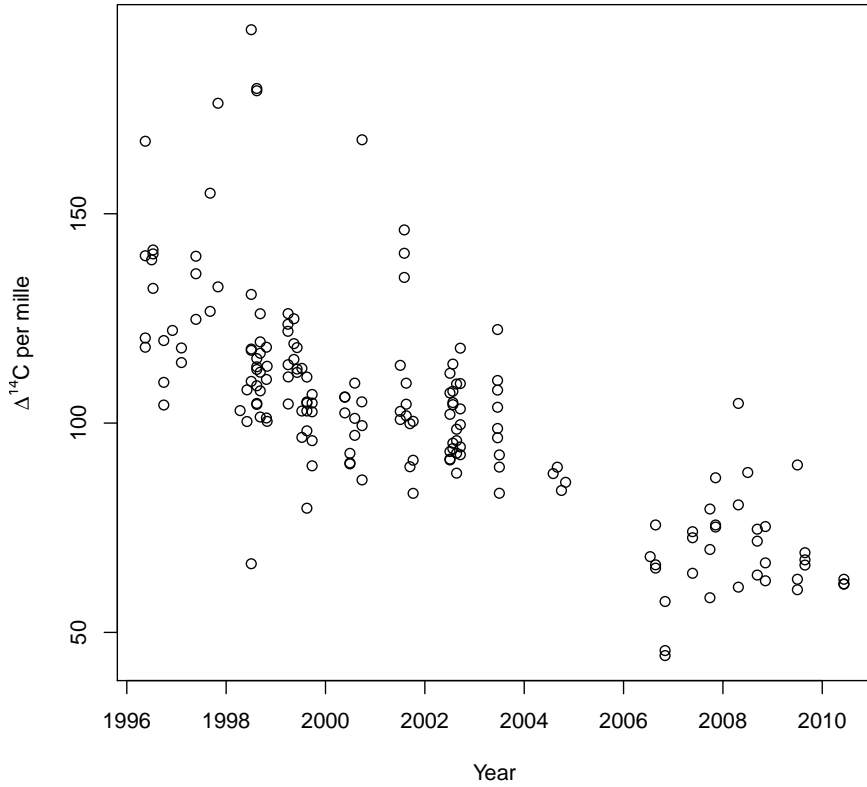


Figure 2: $\Delta^{14}\text{C}$ value of the respired CO₂ in a temperate broadleaf forest at Harvard Forest, USA.

We are interested in fitting the following three-pool model to the data:

$$\frac{d\mathbf{C}(t)}{dt} = I \begin{pmatrix} \gamma_1 \\ \gamma_2 \\ 0 \end{pmatrix} + \begin{pmatrix} -k_1 & 0 & 0 \\ a_{21} & -k_2 & 0 \\ 0 & 0 & -k_3 \end{pmatrix} \begin{pmatrix} C_1 \\ C_2 \\ C_3 \end{pmatrix}. \quad (1)$$

where γ_1 and γ_2 are known. However the parameter a_{21} is not independent from k_2 , which is a condition for the parameter estimation with FME. We therefore formulate the model again based on the independent flux rates.

The radiocarbon content of CO₂ in the atmosphere is necessary for running the model, because it informs us about the concentration and rate of radiocarbon input to the soil. For this example, we will use the dataset `C14Atm_NH` provided with `SoilR`, but other provided datasets such as `Hua2013` can also be used.

First, we define the points in time to run the model from the atmospheric radiocarbon dataset

```
time=C14Atm_NH$YEAR
t_start=min(time)
t_end=max(time)
```

To create the vector of input fluxes we need to create a new object of class `InFluxes`. For our particular model, input fluxes to the C_1 and C_2 pools are created by this command

```
inputFluxes=InFluxes( c(270,150,0))
```

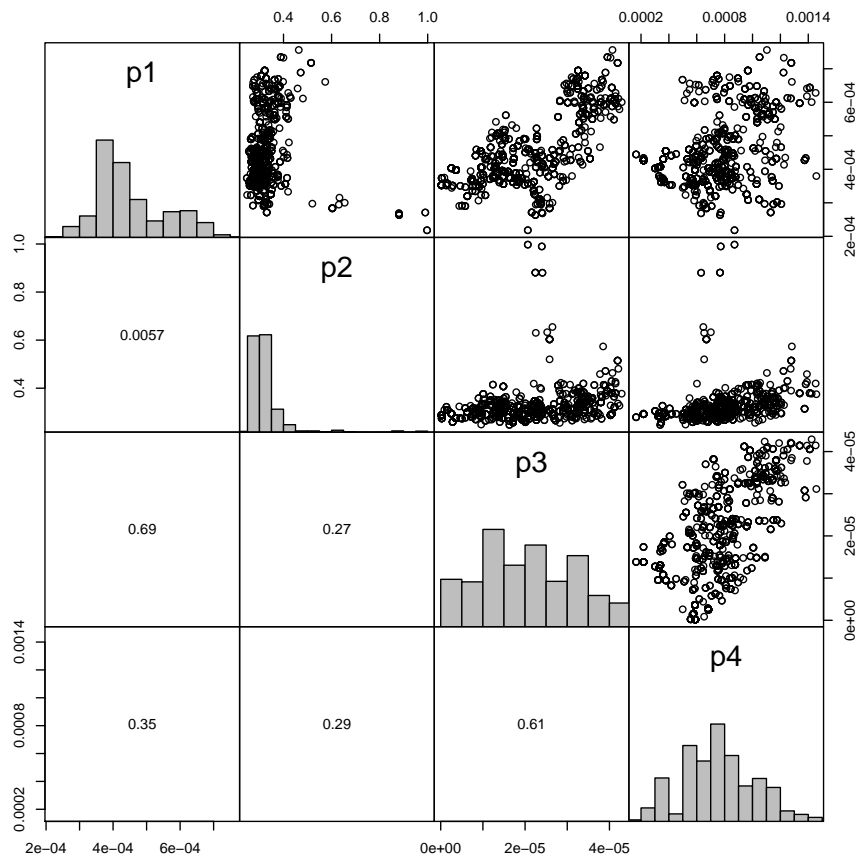
assuming that pool 1 receives 270 gC m² yr⁻¹ and pool 2 150 gC m² yr⁻¹. The initial amount of carbon is created by aggregating the organic and mineral pools for this site reported in Sierra et al. (2012)

```
C0=c(390,220+390+1376,90+1800+560)
```

We can now create the complete `Model` object in `SoilR`. As before we wrap it in a function that takes as arguments a set of parameters and returns the $\Delta^{14}\text{C}$ value of the respired carbon for the desired times.

```
Fc=BoundFc(C14Atm_NH,lag=0,format="Delta14C")
Mod1<-function(ks,pass=TRUE){
  At=ConstLinDecompOp(
    internal_flux_rates=ConstantInternalFluxRateList_by_PoolIndex(
      list(
        "1_to_2"=ks[[4]]
        #, "1_to_3"=ks[[5]]
      )
    )
    ,out_flux_rates=ConstantOutFluxRateList_by_PoolIndex(
      list(
        "1"=ks[[1]]
        ,"2"=ks[[2]]
        ,"3"=ks[[3]]
      )
    )
    ,numberOfPools = 3
  )
  mod=GeneralModel_14(
```

[illegible]



```
summary(Fit)
```

```
##
## Parameters:
##      Estimate Std. Error   t value Pr(>|t|)
## [1,] 2.182e-04  9.378e-05    2.327   0.0212 *
## [2,] 9.973e-01  3.114e-05 32027.873 < 2e-16 ***
## [3,] 2.065e-05  1.504e-03    0.014   0.9891
## [4,] 8.707e-04  1.827e-04    4.765 4.11e-06 ***
## ---
## Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
##
## Residual standard error: 21.15 on 166 degrees of freedom
##
## Parameter correlation:
##      [,1]  [,2]  [,3]  [,4]
## [1,]  1.0000  0.2635 -0.3812 -0.7292
## [2,]  0.2635  1.0000 -0.5405 -0.5118
## [3,] -0.3812 -0.5405  1.0000  0.6037
## [4,] -0.7292 -0.5118  0.6037  1.0000
```


Acknowledgements

We would like to thank Saadat Malghani for performing the laboratory incubation study. Susan E. Trumbore provided the radiocarbon data for Harvard Forest and gave important support and insights for the development of this project. Funding from the Max Planck Society.

References

- Schädel, C., Y. Luo, R. David Evans, S. Fei, and S. Schaeffer. 2013 Separating soil CO₂ efflux into C-pool-specific decay rates via inverse analysis of soil incubation data. *Oecologia* 171: 721–732.
- Sierra, C.A., Saadatullah Malghani and M. Müller. 2015 Model structure and parameter identification of soil organic matter models *Soil Biology and Biochemistry*, 90: 197-203.
- Sierra, C.A., M. Müller, and S. E. Trumbore. 2012. Models of soil organic matter decomposition: the SoilR package, version 1.0. *Geosci. Model Dev.*, 5: 1045–1060.
- Sierra, C.A., S. E. Trumbore, E. A. Davidson, S. D. Frey, K. E. Savage, and F. M. Hopkins. 2012. Predicting decadal trends and transient responses of radiocarbon storage and fluxes in a temperate forest soil. *Biogeosciences*, 9: 3013–3028.
- Soetaert K. and T. Petzoldt. 2010. Inverse modelling, sensitivity and monte carlo analysis in R using package FME. *Journal of Statistical Software*, 33: 1–28.