

Bodenatmung im Nationalpark Hainich

Erstellung eines statistischen Modells und Diskussion von Fehlentscheidungen im Rahmen der Variablen Selektion

Sichtling C

Loos D

Jajali R

14. August 2016

PROJEKTARBEIT

im Rahmen des Moduls *Statistische Verfahren*
an der Friedrich-Schiller-Universität Jena

Die Freisetzung von CO_2 in die Atmosphäre ist ein fundamentaler Bestandteil von Modellen zur Berechnung klimatischer Phänomene. Die *Bodenatmung* ist hierbei der dominierende Prozess auf Landgebieten. Zur Modellierung dieses Prozesses wurde im Nationalpark Hainich Messdaten im Einfluss von ca. 30 Variablen erhoben. Daraus wird nun ein statistisches Modell erstellt. Bei der Variablen Selektion können statistische Unsicherheiten zu Fehlern führen. Dies wird im folgendem diskutiert.

Inhaltsverzeichnis

1	Einleitung	3
1.1	Bodenatmung als geophysikalischer Prozess	3
1.2	Statistische Grundlagen	3
1.2.1	Korrelation	3
1.2.2	F-Statistik	4
2	Erstellung eines Modells zur Bodenatmung	5
2.1	Korrelationsanalyse	5
2.2	Transformationen	5
2.3	Variablen Selektion	5

2.4	Ergebnis	5
2.5	Umsetzung mit R	5

1 Einleitung

1.1 Bodenatmung als geophysikalischer Prozess

[1]

1.2 Statistische Grundlagen

1.2.1 Korrelation

Unter *Korrelation* versteht man die „Wechselbeziehung“ zweier Zufallsgrößen. Dieser Zusammenhang kann entweder *linear* (Korrelationskoeffizient nach Pearson) oder lediglich *monoton* sein. Ist eine Korrelation nicht gegeben, so scheint die Zufallsgröße als ungeeigneter Prädiktor für die jeweils andere Variabel. Kausale Beziehungen erfordern Korrelation.

Der *lineare, empirische Korrelationskoeffizient nach Person* zwischen den Variablen X und Y wird definiert durch:

$$Kor_e(x, y) := \varrho_e(x, y) := r_{xy} := \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2 \cdot \sum_{i=1}^n (y_i - \bar{y})^2}} \quad (1)$$

mit den empirischen Mittelwert \bar{x} aus n Messwerten von X . Möchte man lediglich ein *lineares* Modell erstellen, so ist die Pearson-Korrelation zu wählen.

Eine allgemeinere Beschreibung der Korrelation ist die nach *Spearsman*, welche auf Ranglisten beruht:

$$r_s = \frac{\sum_i (rg(x_i) - \bar{rg}_x)(rg(y_i) - \bar{rg}_y)}{\sqrt{\sum_i (rg(x_i) - \bar{rg}_x)^2} \sqrt{\sum_i (rg(y_i) - \bar{rg}_y)^2}} \quad (2)$$

$$= \frac{\frac{1}{n} \sum_i (rg(x_i)rg(y_i)) - \bar{rg}_x \bar{rg}_y}{s_{rg_x} s_{rg_y}} \quad (3)$$

$$= \frac{Cov(rg_x, rg_y)}{s_{rg_x} s_{rg_y}} \quad (4)$$

Somit lassen sich auch nicht-lineare Korrelationen in Datensätzen erkennen.

Insbesondere im Betrachtung des exponentiellen Verhaltens in Abhängigkeit von der Temperatur ist diese Art der Korrelation wichtig. Zur Berechnung des Anteils an der *erklärten* Varianz in linearen Modellen allerdings kann diese Variante nicht verwendet werden.

1.2.2 F-Statistik

Im Rahmen der Variablenselektion bei der Erstellung des Modells wird immer wieder die Hypothese überprüft, ob ein erweitertes Modell *signifikant* besser ist:

$$F = \frac{\frac{RSS_1 - RSS_2}{p_2 - p_1}}{\frac{RSS_2}{n - p_2}} \quad (5)$$

mit den quadrierten Residuen RSS , den Anzahl an Features p_i des ursprünglichen Modells M_1 und das um ein weiteres Feature erweiterte Modell M_2 .

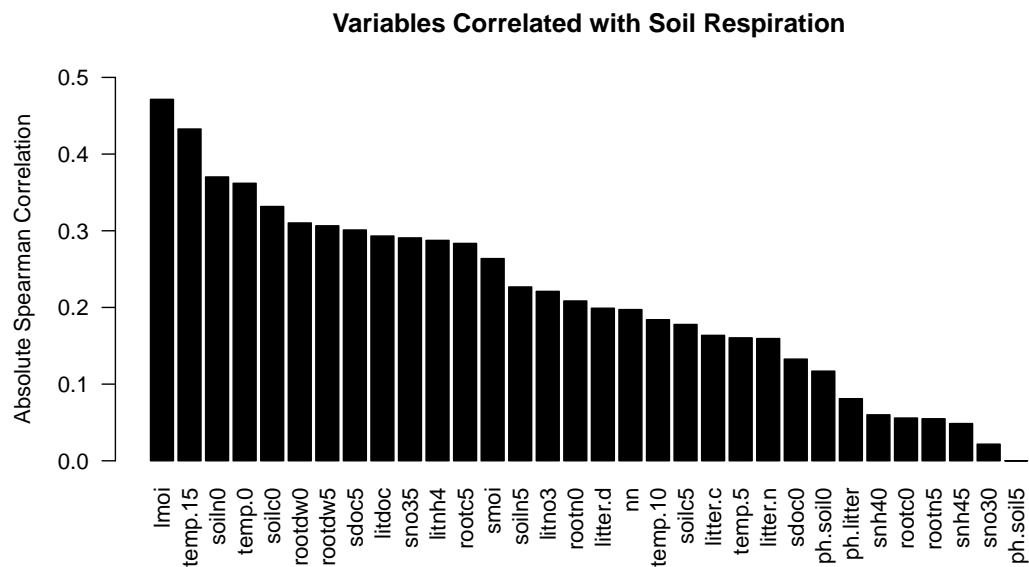


Abbildung 1: **Spearsman-Korrelation** der Einflussgrößen mit der Bodenatmung.

2 Erstellung eines Modells zur Bodenatmung

2.1 Korrelationsanalyse

2.2 Transformationen

$\exp(Tmp)$ anstatt Tmp ?

2.3 Variabelnselektion

2.4 Ergebnis

Qualität des gewählten Modells

2.5 Umsetzung mit R

```

1 ## correlation
2 #correlation from soil.res with all others
3 hainich.r <- abs(cor(hainich, method = "pearson"))["soil.res",-1]
4 hainich.r.ordered <- hainich.r[order(hainich.r, decreasing = T)]
5 barplot(hainich.r.ordered, las = 2, ylim = c(0,0.5), col = "black"
  ,

```

```
6 ylab = "Absolute_Pearson_Correlation", main = "Variables_
  Correlated_with_Soil_Respiration")
7
8 ## scatter plot
9 pairs(~ rootdw0 + smoi + nn,
10 data = hainich, main = "Scatterplot_Matrix")
```

Literatur

[1] N. Jr. My article, 2006.