

# Analyses for 2<sup>nd</sup> draft

Cape vs SWA

*Ruan van Mazijk*

2019-06-22

## Preamble/outline

Here I layout the “new”, second incarnation of the analyses as discussed over the course of May/June 2019, following the first draft of the manuscript.

To reiterate that manuscript, we hypothesise that the greater vascular plant species richness of the GCFR compared to that of the SWAFR is explained by the regions’ difference in environmental heterogeneity.

The proposed “story” of questions for the analyses is as follows:

1. Is the GCFR more heterogeneous environmentally than the SWAFR, and does the scale of that heterogeneity differ to that of the SWAFR?
2. Do the regions differ w.r.t. the species richness of both HDS and QDS cells, and, for HDS cells’ richness ( $S_{HDS}$ ), does the explanatory power of mean QDS richness ( $S_{QDS}$ ) and turnover ( $T_{QDS}$ ) differ between the regions?
3. Does heterogeneity explain differences in richness and turnover between the regions?

## 1. Environmental heterogeneity & scale

Is the GCFR more heterogeneous environmentally than the SWAFR, and does the scale of that heterogeneity differ to that of the SWAFR?

In order to determine which region is more environmentally heterogeneous, and what scales heterogeneity is most pronounced, we calculated a measure of environmental heterogeneity at various spatial scales (namely: the base data resolution ( $0.05^\circ \times 0.05^\circ$ ), eighth- (EDS), quarter- (QDS), half- (HDS) and three-quarter-degree-squares (3QDS)).

Environmental “roughness” in both regions was calculated, in moving  $3 \times 3$  cell windows, as the average absolute difference between cells and their (usually) 8 neighbours. Alternatively, for a focal cell  $x^*$ , the roughness is based on  $x_1, x_2, \dots, x_i, \dots, x_8$  neighbour cells as:

$$Roughness(x^*) = f \begin{pmatrix} x_1 & x_2 & x_3 \\ x_4 & x^* & x_5 \\ x_6 & x_7 & x_8 \end{pmatrix} = \frac{1}{8} \sum_i |x^* - x_i|$$

In R, this is implemented this as follows:

```
roughness <- function(x) {  
  raster::focal(x, matrix(1, nrow = 3, ncol = 3), function(x) {  
    focal_cell <- x[5]  
    focal_exists <- (!is.na(focal_cell)) & (!is.nan(focal_cell))  
    if (focal_exists) {  
      neighbour_exists <- (!is.na(x)) & (!is.nan(x)) & (x != focal_cell)  
      neighbour_cells <- x[neighbour_exists]  
      return(mean(abs(focal_cell - neighbour_cells)))  
    } else {
```

```

    return(NA)
  }
})
}

```

Following this, the various forms environmental heterogeneity were ordinated using principal component analysis (PCA), to summarise a major axis of heterogeneity in each region (Figure 1). Portions of the data matrices for each scale for these PCAs are shown in Table 1.

Both the actual environmental heterogeneity values and the principal component of heterogeneity were then compared between the GCFR and SWAFR using common language effect sizes (*CLES*). The *CLES* of GCFR vs SWAFR heterogeneity values was regressed against the spatial scale at which it was calculated using simple linear regression (Figure 2, Table 2).

We can see that PDQ, NDVI, pH and, arguably, elevation are all consistently more heterogeneous in the GCFR than in the SWAFR, regardless of spatial scale (Figure 2). The GCFR is more heterogeneous at finer scales in terms of MAP, surface temperature, CEC and soil carbon (Figure 2). Notably, the GCFR is more pronouncedly heterogeneous at broad scales in terms of clay (Figure 2). In general (i.e. regarding PC1; Figure 2), the GCFR is more environmentally heterogeneous than the SWAFR, and particularly so at fine spatial scales.

Table 1: Portions of the data matrices used in the PCA for this section of the analysis, where roughness values were  $\log(x + 1)$ -transformed to ensure normality.

region	Elevation	MAP	PDQ	Surface.T	NDVI	CEC	Clay	Soil.C	pH
GCFR	5.19	2.52	0.72	1.32	15.13	1.14	1.2	2.46	1.36
GCFR	5	2.7	0.61	1.16	15.01	1.11	1.11	1.74	1.83
GCFR	4.86	2.55	0.72	1.17	15.08	1.18	1.4	1.79	1.65
...	...	...	...	...	...	...	...	...	...
SWAFR	3.27	2.77	1.1	0.71	14.91	0.31	1.19	1.59	0.48
SWAFR	2.36	2.41	1.15	0.7	14.28	0.67	1.29	2.03	1.3
SWAFR	2.86	1.98	1.17	1.09	13.58	0.73	2.27	2.4	2.58

Table 2: Slopes and associated *P*-values from simple linear regressions of *CLES* against scale for each form of environmental roughness (Figure 2).

Variable	Slope	<i>P</i>	
Elevation	0.044	0.016	*
MAP	-0.313	0.020	*
PDQ	0.010	0.387	
Surface.T	-0.330	0.026	*
NDVI	0.032	0.459	
CEC	-0.126	0.063	.
Clay	0.243	0.013	*
Soil.C	-0.298	0.003	*
pH	-0.010	0.756	
PC1	-0.172	0.010	*

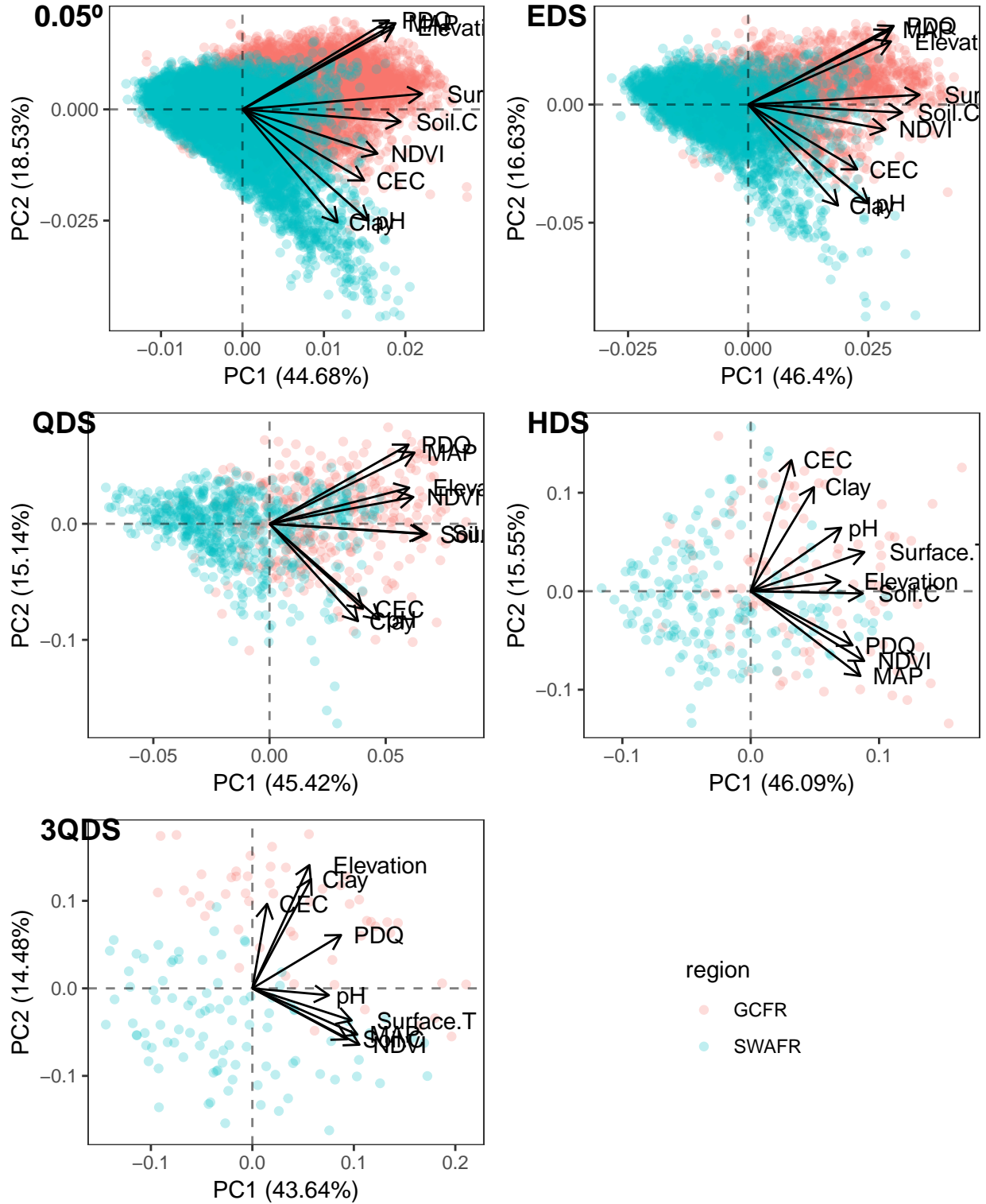


Figure 1: Scatter plots of the first and second principal components (PC1, PC2) of environmental heterogeneity following principal components analyses (PCAs) of the various forms of environmental heterogeneity, repeated at the five spatial scales. The proportion of variation accounted for by each axis is denoted in parentheses. Arrows (labelled) denote the rotational loading of a given form of environmental heterogeneity. Note, the signs of loadings on PC1 have been forced to be positive, while the signs of loadings on PC2 are arbitrary.

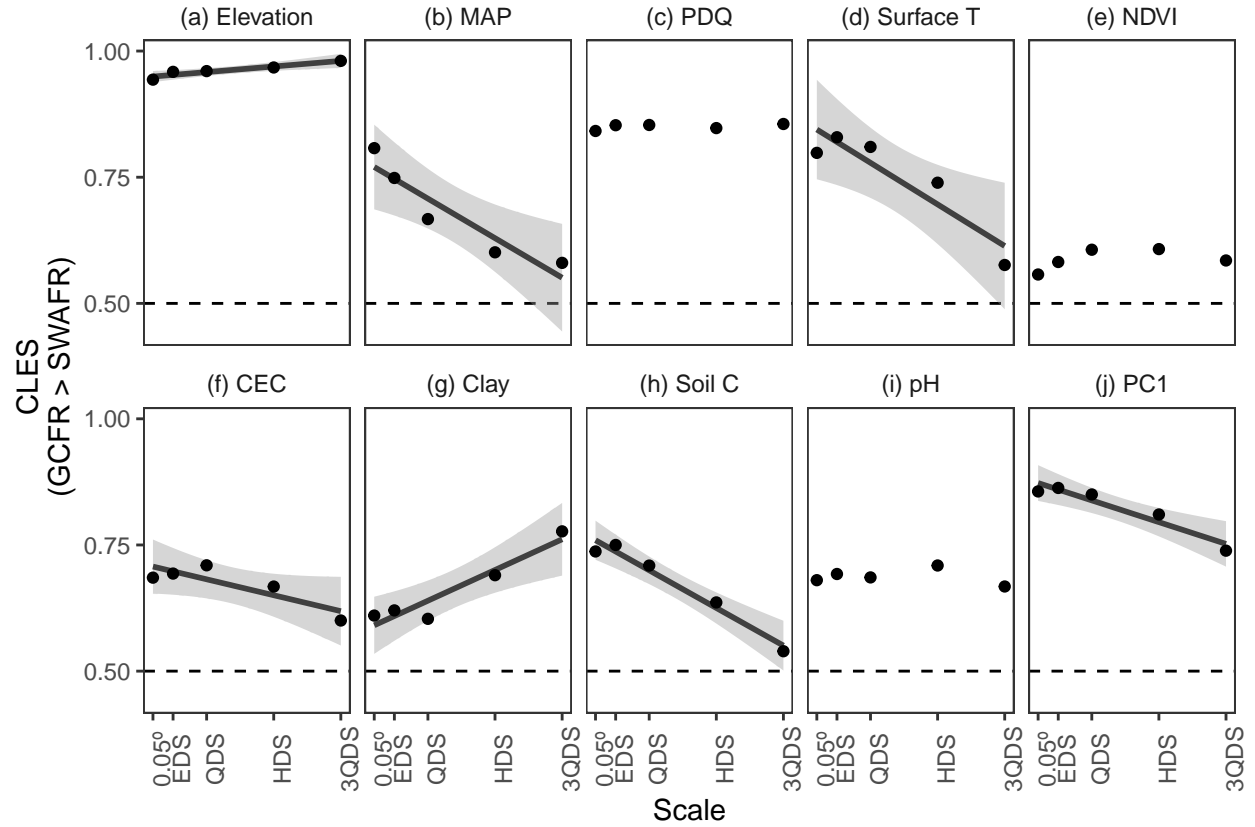


Figure 2: Simple linear regressions of the common language effect size ( $CLES$ ) of various forms of environmental heterogeneity (a–i), and the first principal component of heterogeneity (j; see Figure 1), where the  $CLES$  is treated as the effect of GCFR relative to SWAFR values. Only significant or marginally significant fits are plotted (Table 2). Grey bands denote 95% confidence intervals about the fitted lines. Across spatial scales, all  $CLES$  values differed significantly from zero following two-sided  $t$ -tests ( $P < 0.001$ ).

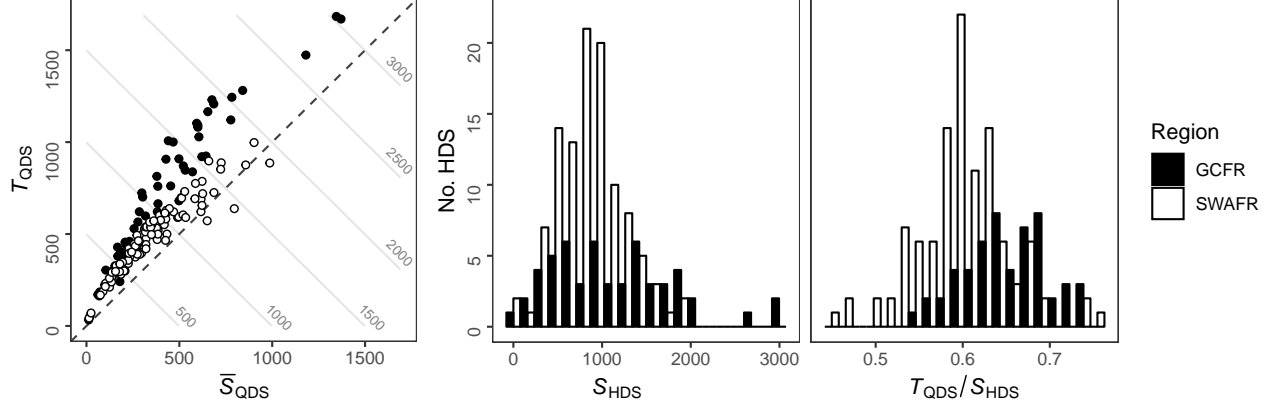


Figure 3: (a) Scatter plot of mean QDS-scale richness ( $\bar{S}_{QDS}$ ) and turnover ( $T_{QDS}$ ) with contour lines denoting the  $S_{HDS}$  that would arise as their sum (i.e. increasing from lower-left to upper-right). Distributions of (a) HDS-scale species richness ( $S_{HDS}$ ) and (b) the turnover partition of that richness ( $T_{QDS}/S_{HDS}$ ).

## 2. Species richness & turnover

Do the regions differ w.r.t. the species richness of both HDS and QDS cells, and, for HDS cells' richness ( $S_{HDS}$ ), does the explanatory power of mean QDS richness ( $S_{QDS}$ ) and turnover ( $T_{QDS}$ ) differ between the regions?

To tackle this question, I compare measures of species richness and turnover between the regions. Species richness at the HDS-scale ( $S_{HDS}$ ) can be partitioned into the average richness of the constituent QDS in HDS ( $\bar{S}_{QDS}$ ) and species turnover ( $T_{QDS}$ ) defined<sup>1</sup> as:

$$T_{QDS} = S_{HDS} - \bar{S}_{QDS}$$

The distributions of these data are presented in Figure 3. To test for significant differences between GCFR and SWAFR values, I use Mann-Whitney  $U$ -tests and  $CLES$  (Table 3), as most of the variables deviate significantly from normality (Shapiro-Wilk normality test;  $P < 0.05$ ).

Additionally, a visualisation of how  $S_{HDS}$  is partitioned into  $\bar{S}_{QDS}$  and  $T_{QDS}$  is presented in Figure 4.

We can conclude that broad scale species richness (i.e. that at the HDS scale) is more strongly driven by turnover between areas (i.e. QDS) than so in the SWAFR.

Table 3: Results of Mann-Whitney  $U$ -tests and the  $CLES$  of GCFR vs SWAFR for various species richness and turnover metrics.

Metric	$CLES$	$P_U$
$S_{HDS}$	0.612	0.020
$S_{QDS}$	0.595	< 0.001
$T_{QDS}/S_{HDS}$	0.784	< 0.001

## 3. Relating heterogeneity to species richness & turnover

Does heterogeneity explain differences in richness and turnover between the regions?

<sup>1</sup>following Whittaker's original additive definition:  $\gamma = \alpha + \beta$

Here I fit various linear regressions of richness and turnover as functions of environmental heterogeneity across the two regions. The richness and turnover measures used are the same as in the previous section, while the environmental heterogeneity was recalculated in the same grid-wise fashion as the richness and turnover measures. These analyses were carried out at both the HDS- and QDS-scales, insofar as species occurrence data from GBIF is only accurate to the QDS-scale. These analyses were only carried out on HDS-scale data for HDS-cells that contained four QDS-cells, and similarly for QDS-scale data for QDS-cells that contained four EDS-cells.

Environmental “roughness” here was calculated for each HDS- and QDS-cell in both regions as the mean of each consituent QDS- and EDS-cell’s mean absolute difference in environmental conditions from the other three cells within that HDS- or QDS-cell.

In other words, roughness was calculated by first calculating the average absolute-difference in environmental values between each QDS and it’s three neighbours in a given HDS. Then, these four values (assuming four QDS in an HDS) are averaged. This roughness index is presented mathematically below. This index allows each of the four values to be similarly independent, and thus more sutiable for our averaging and analyses, as opposed to if it were simly the direct average of pairwise differences [expand?].

$$Roughness_{cellular}(\{x_1, x_2, x_3, x_4\}) = \frac{1}{4} \sum_i f(x_i) = \frac{1}{4} \sum_i \left( \frac{1}{3} \sum_{j \neq i} |x_i - x_j| \right)$$

In R, this is implemented this as follows:

```
roughness_cellular <- function(x) {
  out <- vector(mode = "numeric", length = length(x))
  for (i in seq_along(x)) {
    out[[i]] <- mean(abs(x[i] - x[-i]))
  }
  mean(out)
}
```

The models I fit here are as follows:

Response	Variables	Region term
$S_{HDS}$	Separate	Separate
"	"	None
"	"	Additive
"	"	Interaction
"	AIC-set	Separate
"	"	None
"	"	Additive
"	"	Interaction
"	PC1	Separate
"	"	None
"	"	Additive
"	"	Interaction
$S_{QDS}$	Separate	Separate
"	"	None
"	"	Additive
"	"	Interaction
"	AIC-set	Separate
"	"	None
"	"	Additive

Response	Variables	Region term
"	"	Interaction
"	PC1	Separate
"	"	None
"	"	Additive
"	"	Interaction
$T_{\text{HDS}}/S_{\text{HDS}}$	Separate	Separate
"	"	None
"	"	Additive
"	"	Interaction
"	AIC-set	Separate
"	"	None
"	"	Additive
"	"	Interaction
"	PC1	Separate
"	"	None
"	"	Additive
"	"	Interaction

### 3.1. Separate environmental variable models

#### 3.1.1. With $S_{\text{HDS}}$ as the response

Table 5: Results of separate simple linear regressions of  $S_{\text{HDS}}$  against environmental heterogeneity variables.

Environmental predictor	Model	$\Delta AIC$		$P_{\text{slope}}$	$P_{\text{region}}$	$P_{\text{slope:region}}$
Elevation	None	0.000	*	< 0.050		
	Additive	1.627		< 0.050	0.546	
	Interaction	0.823		0.181	0.899	0.099
MAP	None	6.950		< 0.050		
	Additive	8.874		< 0.050	0.785	
	Interaction	0.000	*	< 0.050	0.334	< 0.050
PDQ	None	0.000	*	< 0.050		
	Additive	1.944		< 0.050	0.815	
	Interaction	1.737		< 0.050	0.915	0.143
Surface T	None	0.000	*	< 0.050		
	Additive	1.523		< 0.050	0.494	
	Interaction	2.017		0.147	0.368	0.226
NDVI	None	2.491		< 0.050		
	Additive	0.205	*	< 0.050	< 0.050	
	Interaction	0.000		< 0.050	0.090	0.143
CEC	None	7.612		0.141		
	Additive	0.000	*	0.827	< 0.050	
	Interaction	1.644		0.720	< 0.050	0.556
Clay	None	7.972		< 0.050		
	Additive	0.000	*	< 0.050	< 0.050	
	Interaction	1.656		0.478	< 0.050	0.563
Soil C	None	0.000	*	< 0.050		
	Additive	0.508		< 0.050	0.227	
	Interaction	2.439		< 0.050	0.318	0.795

Environmental predictor	Model	$\Delta AIC$		$P_{\text{slope}}$	$P_{\text{region}}$	$P_{\text{slope:region}}$
pH	None	6.320		< 0.050		
	Additive	0.000	*	0.383	< 0.050	
	Interaction	1.427		0.282	< 0.050	0.455

### 3.1.2. With $S_{\text{QDS}}$ as the response

Table 6: Results of separate simple linear regressions of  $S_{\text{QDS}}$  against environmental heterogeneity variables.

Environmental predictor	Model	$\Delta AIC$		$P_{\text{slope}}$	$P_{\text{region}}$	$P_{\text{slope:region}}$
Elevation	None	0.000	*	< 0.050		
	Additive	0.806		< 0.050	0.276	
	Interaction	2.129		< 0.050	0.179	0.412
MAP	None	0.000	*	< 0.050		
	Additive	1.974		< 0.050	0.871	
	Interaction	2.541		< 0.050	0.513	0.233
PDQ	None	12.963		< 0.050		
	Additive	14.523		< 0.050	0.509	
	Interaction	0.000	*	< 0.050	0.810	< 0.050
Surface T	None	0.000	*	< 0.050		
	Additive	0.274		< 0.050	0.190	
	Interaction	1.528		< 0.050	0.134	0.390
NDVI	None	24.274		< 0.050		
	Additive	12.991		< 0.050	< 0.050	
	Interaction	0.000	*	< 0.050	< 0.050	< 0.050
CEC	None	23.788		0.188		
	Additive	0.000	*	0.733	< 0.050	
	Interaction	1.907		0.659	< 0.050	0.761
Clay	None	22.889		< 0.050		
	Additive	0.204	*	0.207	< 0.050	
	Interaction	0.000		0.537	< 0.050	0.139
Soil C	None	9.586		< 0.050		
	Additive	5.067		< 0.050	< 0.050	
	Interaction	0.000	*	< 0.050	0.213	< 0.050
pH	None	24.721		0.362		
	Additive	0.000	*	0.703	< 0.050	
	Interaction	1.899		0.909	< 0.050	0.752

### 3.2.3. With $T_{\text{QDS}}/S_{\text{HDS}}$ as the response

Table 7: Results of separate simple linear regressions of  $T_{\text{QDS}}/S_{\text{HDS}}$  against environmental heterogeneity variables.

Environmental predictor	Model	$\Delta AIC$		$P_{\text{slope}}$	$P_{\text{region}}$	$P_{\text{slope:region}}$
Elevation	None	53.703		< 0.050		
	Additive	12.519		< 0.050	< 0.050	
	Interaction	0.000	*	0.926	< 0.050	< 0.050
MAP	None	55.720		0.276		
	Additive	1.382	*	< 0.050	< 0.050	



Environmental predictor	Model	$\Delta AIC$		$P_{\text{slope}}$	$P_{\text{region}}$	$P_{\text{slope:region}}$
PDQ	Interaction	0.000		0.359	< 0.050	0.070
	None	42.488		< 0.050		
	Additive	8.905		< 0.050	< 0.050	
Surface T	Interaction	0.000	*	0.981	< 0.050	< 0.050
	None	36.444		< 0.050		
	Additive	8.396		0.186	< 0.050	
NDVI	Interaction	0.000	*	0.205	< 0.050	< 0.050
	None	38.893		0.869		
	Additive	0.000	*	< 0.050	< 0.050	
CEC	Interaction	1.104		0.605	< 0.050	0.351
	None	20.735		< 0.050		
	Additive	0.000	*	0.237	< 0.050	
Clay	Interaction	1.995		0.535	< 0.050	0.946
	None	39.461		0.070		
	Additive	0.000	*	< 0.050	< 0.050	
Soil C	Interaction	1.950		0.129	< 0.050	0.826
	None	51.192		0.691		
	Additive	8.043		< 0.050	< 0.050	
pH	Interaction	0.000	*	0.283	< 0.050	< 0.050
	None	24.761		< 0.050		
	Additive	0.000	*	0.168	< 0.050	
	Interaction	1.994		0.542	< 0.050	0.938

### 3.x. Separate-regions models with combinations of variables

Table 8: Results of bi-directional stepwise multiple linear regressions of three richness and turnover responses in the against additive combinations of environmental heterogeneity variables. The step-wise regression procedure started with all variables included. (See Figure 5 for a graphical representation.)

Region	Response	Predictor	Slope	$P_{\text{slope}}$	
GCFR	$S_{\text{HDS}}$	Clay	185.456	0.019	*
GCFR	$S_{\text{HDS}}$	MAP	738.358	0.000	*
GCFR	$S_{\text{HDS}}$	pH	-322.625	0.006	*
GCFR	$S_{\text{QDS}}$	MAP	136.688	0.003	*
GCFR	$S_{\text{QDS}}$	NDVI	139.568	0.000	*
GCFR	$S_{\text{QDS}}$	PDQ	-45.541	0.147	
GCFR	$S_{\text{QDS}}$	pH	-164.670	0.000	*
GCFR	$S_{\text{QDS}}$	Soil.C	97.764	0.009	*
GCFR	$T_{\text{QDS}}/S_{\text{HDS}}$	Clay	-0.017	0.016	*
GCFR	$T_{\text{QDS}}/S_{\text{HDS}}$	MAP	-0.026	0.010	*
GCFR	$T_{\text{QDS}}/S_{\text{HDS}}$	Soil.C	0.024	0.015	*
SWAFR	$S_{\text{HDS}}$	CEC	-111.775	0.000	*
SWAFR	$S_{\text{HDS}}$	Clay	56.676	0.036	*
SWAFR	$S_{\text{HDS}}$	Elevation	200.297	0.000	*
SWAFR	$S_{\text{HDS}}$	MAP	108.435	0.001	*
SWAFR	$S_{\text{HDS}}$	PDQ	180.511	0.001	*
SWAFR	$S_{\text{HDS}}$	Surface.T	99.867	0.027	*
SWAFR	$S_{\text{QDS}}$	CEC	-28.862	0.012	*

Region	Response	Predictor	Slope	$P_{slope}$	
SWAFR	$S_{QDS}$	Clay	18.683	0.094	.
SWAFR	$S_{QDS}$	Elevation	42.177	0.014	*
SWAFR	$S_{QDS}$	MAP	97.709	0.000	*
SWAFR	$S_{QDS}$	PDQ	116.652	0.000	*
SWAFR	$S_{QDS}$	Surface.T	47.573	0.002	*
SWAFR	$T_{QDS}/S_{HDS}$	CEC	0.014	0.008	*
SWAFR	$T_{QDS}/S_{HDS}$	Clay	-0.011	0.022	*
SWAFR	$T_{QDS}/S_{HDS}$	Elevation	-0.035	0.000	*
SWAFR	$T_{QDS}/S_{HDS}$	MAP	-0.009	0.066	.
SWAFR	$T_{QDS}/S_{HDS}$	PDQ	-0.015	0.113	
SWAFR	$T_{QDS}/S_{HDS}$	pH	0.011	0.020	*
SWAFR	$T_{QDS}/S_{HDS}$	Soil.C	-0.012	0.046	*

Table 9: Adjusted  $R^2$ -values of the models in Table 5.

Response	GCFR $R^2_{adj.}$	SWAFR $R^2_{adj.}$
$S_{HDS}$	0.429	0.510
$S_{QDS}$	0.262	0.323
$T_{QDS}/S_{HDS}$	0.139	0.424

## 3.2. Combined-regions models with individual variables

### 3.2.2. PC1 models

Here, I present my findings with raw R-code, because I don't have the time to format it neatly.

```
m1 <- lm(HDS_richness ~ PC1, HDS)
m2 <- lm(HDS_richness ~ PC1 + region, HDS)
m3 <- lm(HDS_richness ~ PC1 * region, HDS)
my_AIC_table(m1, m2, m3, caption = "Richness (HDS)")
```

```
##      model      AIC delta_AIC w_Akaike
## 1 No region 2433.144    0.000    0.451
## 2 Add. region 2433.558    0.414    0.367
## 3 Int. region 2434.958    1.814    0.182
```

```
# Therefore, "choose" m1 ("no region" model)
```

```
m1 <- lm(QDS_richness ~ PC1, QDS)
m2 <- lm(QDS_richness ~ PC1 + region, QDS)
m3 <- lm(QDS_richness ~ PC1 * region, QDS)
my_AIC_table(m1, m2, m3, caption = "Richness (QDS)")
```

```
##      model      AIC delta_AIC w_Akaike
## 1 No region 9205.960    0.999    0.262
## 2 Add. region 9204.961    0.000    0.432
## 3 Int. region 9205.652    0.691    0.306
```

```
# Therefore, "choose" m1 ("no region" model) (?)
```

```
m1 <- lm(add_turnover ~ PC1, HDS)
m2 <- lm(add_turnover ~ PC1 + region, HDS)
```

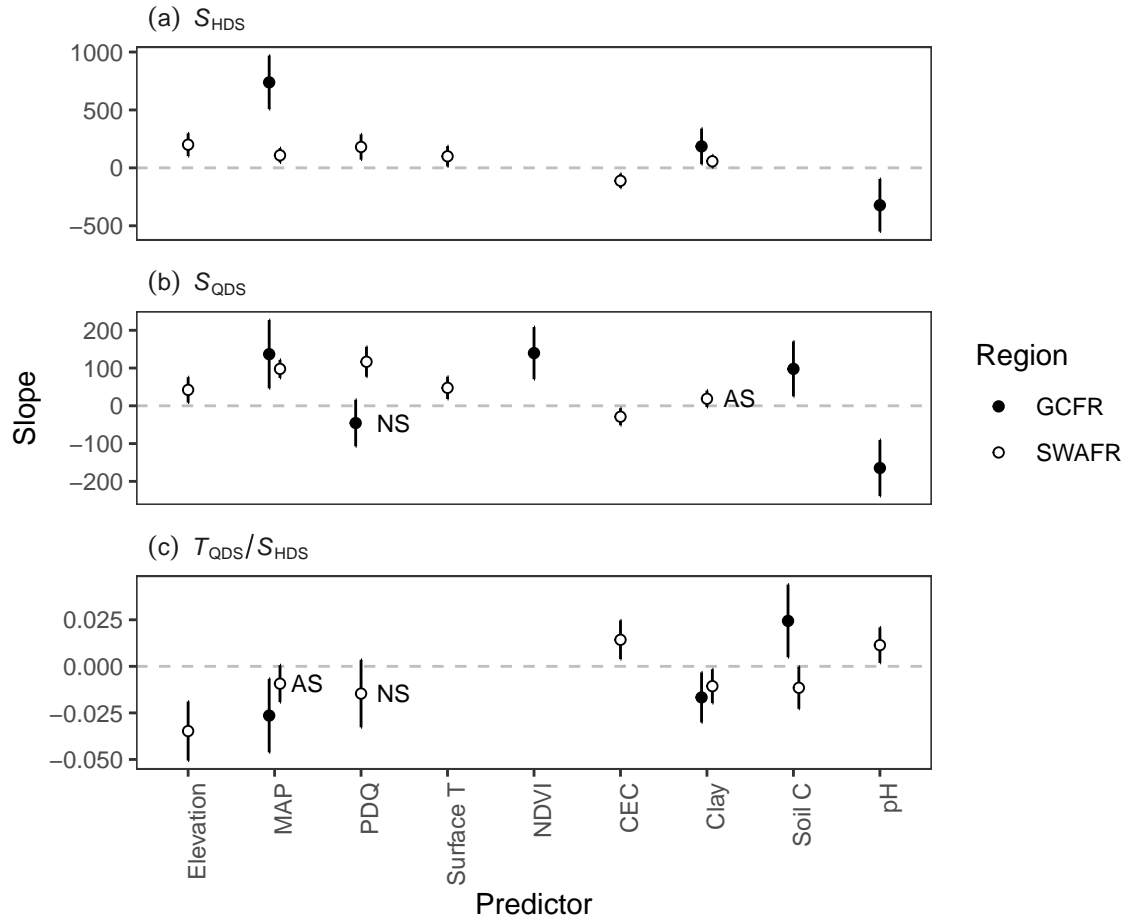


Figure 4: Slopes from Table 5, with error bars denoting 95% confidence intervals about each slope estimate.

```
m3 <- lm(add_turnover ~ PC1 * region, HDS)
my_AIC_table(m1, m2, m3, caption = "Turnover")
```

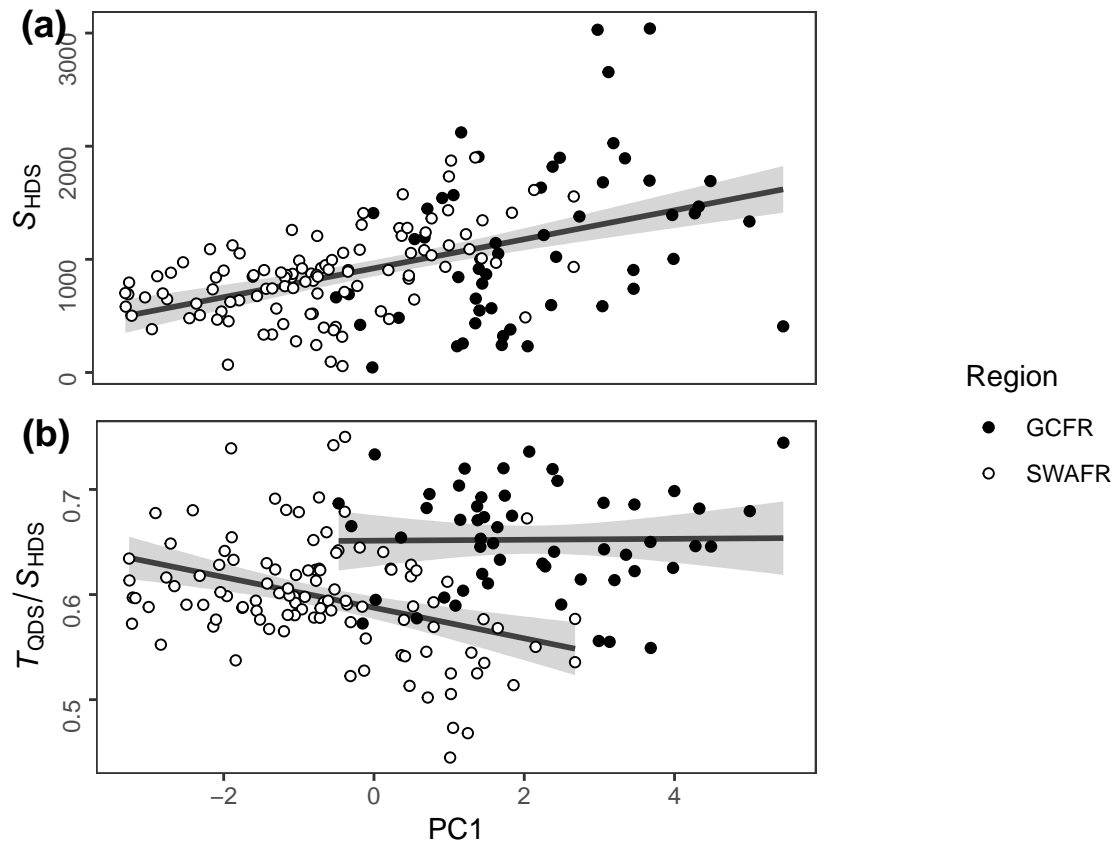
```
##          model      AIC delta_AIC w_Akaike
## 1  No region 2240.186    0.000   0.536
## 2 Add. region 2242.178    1.993   0.198
## 3 Int. region 2241.587    1.401   0.266
```

```
# Therefore, "choose" m1 ("no region" model)
```

```
m1 <- lm(add_turnover_prop ~ PC1 ,      HDS)
m2 <- lm(add_turnover_prop ~ PC1 + region, HDS)
m3 <- lm(add_turnover_prop ~ PC1 * region, HDS)
my_AIC_table(m1, m2, m3, caption = "Turnover (proportional)")
```

```
##          model      AIC delta_AIC w_Akaike
## 1  No region -458.545   45.812   0.000
## 2 Add. region -499.982    4.374   0.101
## 3 Int. region -504.357    0.000   0.899
```

```
# Therefore, "choose" m3 ("int. region" model)
```



### 3.3. Combined-regions models with combinations of variables