

## **Analiza și modelarea statistică multivariată a comunităților; Analiza de ordonare**

Foarte frecvent în ecologie analizăm comunități (sisteme suprapopulaționale), în habitate caracterizate prin diferiți parametri calitativi și cantitativi, în dorința de a analiza modul în care acestea răspund la factorii de mediu. În acest gen de studii dorim să răspundem la diverse întrebări sau probleme, cum ar fi:

- modul în care fiecare populație reacționează (își modifică parametri ecologici cantitativi) la variația graduală a unui factor de mediu (sau mai mulți);
  - modificarea unei populații analizată pe diferiți gradienti ai mediului;
  - răspunsul mai multor populații (specii) la variația graduală a unui număr oarecare de caracteristici ale habitatului;
  - care factori ai mediului explică mai bine și în ce mod se reflectă aceștia în compoziția comunității (cum influențează gradientii mediului comunitatea)?
  - cum se clasifică variabilele de mediu în funcție de gradul de explicare a compoziției speciilor (structura comunității)?
  - cât din variația compoziției specifice (a variabilelor dependente sau de răspuns) este explicată de către variabilele de mediu evaluate?
- ... și multe altele ...

La aceste întrebări putem prea rar răspunde dacă le descompunem în părți sau în elemente constituente; frecvent nici nu putem să adoptăm un model reduționist, în care să analizăm separat fiecare variabilă sau să căutăm modele de regresie între variația fiecărei specii în relație cu modificarea fiecărei variabile de mediu.

Dacă studiem modul în care o variabilă ecologică (abundența unei populații de exemplu) este relaționată de o variabilă de mediu (altitudinea, salinitatea etc.) o analiză de regresie, liniară sau neliniară, ar putea fi suficientă, iar relația s-ar putea reprezenta grafic sub forma unui plan bidimensional (2D), în care ordonata este prima, iar abscisa definește cea de-a doua variabilă (care descrie variația mediului). Când dorim să analizăm relații multidimensionale (nD) între diverse specii și mai multe variabile de mediu (independente sau de habitat), precum și să reprezentăm grafic aceste relații, vom fi constrânși să înțelegem și să aplicăm metode sau modele statistice multivariate, dintre care analizele de ordonare fac parte integrantă.

**Analiza de ordonare** este o metodă de studiu a modalităților prin care comunitățile (ansambluri de populații simpatrice și sincronice care interacționează între ele și cu factorii de mediu) răspund la variația graduală a condițiilor de habitat. Deoarece factorii de mediu se modifică gradual, diversele specii vor răspunde prin variația treptată, continuă, a abundenței sau a altui parametru ecologic. Prin urmare o

analiză de ordonare se poate asemăna cu un studiu de gradient, respectiv se pretează la o analiză a schimbărilor factorilor de mediu care le condiționează pe cele ale parametrilor populațiilor. O altă idee este aceea că factorii de mediu sunt frecvent mai mult sau mai puțin corelați, sau sunt grupați în funcție de gradul de interacțiune, speciile modificându-și parametri cantitativi în relație mai puternică cu unii dintre aceștia și mai puțin cu alții, uneori demonstrând chiar o relativă indiferență la variația unora dintre ei.

Comunitatea se modifică neuniform și treptat, deoarece populațiile din care este alcătuită sunt adaptate în mod diferit la factorii de mediu și își modifică relativ independent parametri cantitativi în funcție de aceștia. Dar, atunci când analizăm aceste schimbări, constatăm adesea că răspunsul grupelor de specii, respectiv a comunității întregi, pot fi prevăzute în manifestarea lor, iar gradientii de-a lungul cărora se înregistrează variațiile pot fi sintetizați sub forma unui număr redus de axe sau de gradienti imaginari, care redau sintetic corelația dintre diverșii factori ai mediului.

Câteodată gradientii sau axele de-a lungul cărora se modifică abundența și compoziția speciilor pot fi recunoscuți într-un peisaj, de exemplu sub forma umidității (trecerea de la un habitat acvatic la unul uscat, terestru) sau a altitudinii. Alteori aceste axe nu pot fi identificate ca atare și atunci vorbim de gradienti de modificare (schimbare) a compoziției speciilor.

Studiul variației comunităților se poate face în multe feluri (există metode diverse), dar atunci când analizăm schimbarea treptată și continuă (graduală sau pe unul sau mai mulți gradienti) metodele de ordonare sunt printre cele mai eficiente și adesea singurele capabile să răspundă la întrebările pe care le formulăm. Expresia grafică a acestor metode poartă numele de **diagrame de ordonare** și acestea implică proiectarea speciilor, a probelor (eșantionelor) și/sau a variabilelor (factorilor) de mediu într-un **spațiu de ordonare**, adică pe un plan construit dintr-un număr redus de axe. Regulile de reprezentare și interpretare vor fi studiate în cele ce urmează, dar esența înțelegerii diagramelor este schimbarea continuă a compoziției comunităților de-a lungul unor gradienti, care sunt reprezentați sub forma unor axe, precum și ideea că proximitatea implică similaritate (ceea ce este mai aproape, este mai asemănător sau între entități sunt relații mai strânse). Reprezentarea și interpretarea graficelor (a diagramelor de ordonare) este numai o mică parte a problemei: cel puțin la fel de importantă este modelarea statistică, în sensul testării diverselor ipoteze (de exemplu semnificația relației speciilor cu diverși factori ai mediului, clasificarea acestora în funcție de importanța lor predictivă etc.), al conceptualizării și al identificării relațiilor multiple între diversele variabile, îndepărtarea variației aleatoare și identificarea

modelelor repetitive, legice, de tip corelativ și de tip cauză-efect, între diversele categorii de fenomene.

De la bun început este bine să înțelegem că metodele descrise aici includ o componentă grafică, descriptivă sau explorativă, dar nu se limitează nici pe departe cu aceasta. Tehnicile expuse includ și permit dezvoltarea unor unelte complexe de testare a ipotezelor multiple de lucru, analize de regresie, de varianță, relații sau dependență de timp, măsurători multiple repetate asupra acelorași obiecte, analiza datelor structurate spațial, analiză ierarhică de varianță (ANOVA și MANOVA) etc.. Aceste tehnici permit cercetătorilor să conceapă studii și să ridice întrebări complexe și frecvent mai adecvate problemelor abordate în studiile moderne de ecologie. Prin urmare: deși sunt metode excelente pentru **analiza exploratorie a datelor**, acestea sunt extrem de utile (uneori de neînlocuit) și în **analiza confirmatorie, adică în studiile ipotetic-deductive, respectiv cele bazate pe design experimental complex**.

Amănuntele și formulele matematice vor fi eludate în cea mai mare parte, pentru a simplifica cât mai mult aspectele mecanismelor și pentru a ne concentra mai ales asupra efectelor analizelor datelor ecologice, utilizând pachetul de programe cunoscut sub denumirea de *Canoco*. Acest produs software este concentrat pe reducerea dimensiunii datelor (ordonare), analiza de regresie și combinarea celor două, tehnică cunoscută sub denumirea de **ordonare constrânsă, condiționată sau ordonare canonică**. Obiectivul general este să ajute cercetătorii să înțeleagă și să identifice modele și structuri în datele multivariate, prin relații între variabile și articole. Numele softului *Canoco* provine de la acronimul **CANOnical Community Ordination** și provine din domeniul ecologiei comunităților. Din acest motiv a adoptat o serie de termeni și denumiri din această disciplină, cum ar fi **probă** pentru un articol sau o linie, **specie** pentru o variabilă de răspuns sau variabilă dependentă, precum și **variabilă de mediu** pentru o variabilă explicativă, independentă sau suplimentară.

Scopul **ordonării** este cel de a identifica **axe ale variabilității maxime în compoziția comunității** (numite **axe de ordonare**), pentru un set de probe, precum și de a vizualiza (utilizând **diagrame de ordonare**) structurile de similaritate pentru probe și specii. Este de așteptat ca axele de ordonare să coincidă cu unele variabile de mediu, iar atunci când aceste variabile sunt măsurate, le vom corela cu axele de ordonare. Scopul **ordonării constrânse sau canonice** este de a găsi variabilitatea în compoziția speciilor care poate fi explicată prin măsurarea variabilelor de mediu. Ordonarea poate fi și **parțial constrânsă (sau parțial canonică)**: în analizele parțiale mai întâi se îndepărtează variabilitatea în compoziția speciilor explicată de **covariate sau covariabile** (variabile de mediu definite ca atare), după care se realizează o ordonare canonică a **variabilității reziduale** sau remanente. Există și analiza de

**ordonare hibridă**, în care primele  $x$  axe ( $x$  este uzual 1, 2 sau 3) de ordonare sunt canonice (constrânse) iar celelalte sunt neconstrânse sau necanonice.

### Tipuri de variabile incluse în analiza multivariată a comunităților

Orice set de date primare (observaționale sau experimentale) este format din:

- **variabile ale comunității, pe care le mai numim (poate cam impropriu) ”de specii”, dependente, sau de răspuns** (așezate pe coloane de obicei);
- **valorile acestora în probe (eșantioane, experiențe) numite și articole** (pe linii de obicei); (în termeni de prezență-absență, densitate, efectiv, abundența relativă, indicele Braun-Blanquet etc.). Valorile asociate cu speciile pot fi cantitative dar pot fi și date de prezență-absență (calitative sau de tip 1/0). Corespunzător modului în care definim și evaluăm speciile se va determina și scara variabilelor dependente sau de răspuns: cantitative, semi-cantitative (de exemplu mijlocul claselor de abundență-dominanță relativă prin scara Braun-Blanquet) sau de tip ”boolean” sau binare (adică 1/0, da/nu, prezență/absență).

Frecvent dispunem de valori (articole) și pentru:

- **variabile de mediu, pe care le mai numim și externe, independente, canonice, explicative, de condiționare, de constrângere etc.**

**Variabilele independente** (în mod asemănător cu cele dependente) **pot fi:**

- **cantitative** (numere reale sau numere naturale etc.), de exemplu altitudinea;
- **semicantitative** (ex. impactul antropic pe o scară de la 0 - absent la 10 - foarte puternic); în analize acestea sunt de obicei tratate similar cu variabilele numerice, cantitative;
- **nominale** (de tip **categorie**, numite și **factori**) - pot fi de tip prezență/absență (1/0 sau **variabile binare**), sau pot fi de tip **variabile artificiale, substitutive** (*dummy variables*) sau de tip **indicator**. Legat de ultima categorie - dacă vrem să codificăm tipul de folosință al unui teren - dacă este fâneată, pajiște, alte folosințe sau este abandonat, în *Canoco 4.x* nu putem defini o variabilă cu 4 categorii (adică cele enunțate) ci trebuie să definim 4 variabile, fiecare cu valoare de 0 pentru fals (sau negație) și 1 pentru adevăr (sau afirmație). Astfel, folosința terenului se va codifica prin variabile, cum ar fi **fâneată (da/nu, 1/0), pajiște (1/0), altele (1/0) și părăsit (1/0)**. În *Canoco 5.x* această problemă s-a rezolvat prin posibilitatea de definire a variabilelor factoriale sau a **factorilor**. Astfel, în cazul de mai sus o variabilă de tip *folosința terenului* va fi codificată (de exemplu sub numele sau **eticheta FOLOS**) și va cuprinde valori alfanumerice (deci litere sau combinații), cum ar fi  $F$  pentru articolele care conțin date din probele sau suprafețele de eșantionare situate în fânețe,  $P$  pentru pajiști,  $A$  pentru alte folosințe și  $P$  pentru teren părăsit.

Variabilele explicative pot fi heterogene: unele sunt cantitative, altele pot fi ordinale sau binare, în același fișier.

În analize (și implicit în *Canoco*) fișierele sunt de tip rectangular, fiind constituite **din tabele de specii (ale comunităților) și tabele de variabile explicative** sau de **mediu**. De obicei o **coloană este o variabilă** (specie sau parametru de mediu), iar o linie este **un articol, sau o probă unitară**, adică valorile tuturor variabilelor care provin dintr-un anumit eșantion individual. Relația de corespondență dintre cele două tabele este dată de articole sau de linii, motiv pentru care la fiecare tabel, prima coloană se referă la numărul curent al probelor, iar ordinea (poziția acestora) trebuie să fie aceeași. Adică: un anumit rând trebuie să conțină atât în tabelul de specii cât și în cel de mediu valorile variabilelor din aceeași probă.

Și acum, ceva **teorie despre metodele aplicate în ecologia comunităților** și locul, respectiv rolul, tehnicilor multivariate de analiză a gradientului, construirea și interpretarea diagramelor de ordonare și unele metode relaționate.

Aceste metode trebuie puse în context; nu sunt apărute din senin și studenții de la ecologie sau biologie au avut de-a face atât cu metode de studiu ale comunităților cât și cu tehnici și metode de analiză multivariată. În timpul studiilor de licență, la laboratoarele de Ecologie generală, Ecologia populațiilor, Ecosisteme: structură și funcții, precum și la Modelarea proceselor și sistemelor ecologice (se presupune că) s-au învățat:

- Definirea, abordarea și natura sistemelor ecologice pe diferite niveluri de organizare (abordarea teoretică; identificarea sistemelor ecologice);
- Parametrii ecologici cantitativi ai populațiilor și ai comunităților;
- Grafice: reprezentarea (cantitativă a) comunităților și a dinamicii acestora;
- Analiza de regresie univariată și multivariată, liniară și neliniară.  
(dacă ați uitat despre ce vorbim, recapitulați temele din volumul de Ecologie Practică, Sîrbu și Benedek, ediția a 3-a, 2012)
- Analiza de corelație;
- Testarea ipotezelor;
- Tehnici și metode speciale aplicate în studiul comunităților de păsări și de mamifere;
- Nișa ecologică: lățimea, preferințele și suprapunerea nișelor (ultima fiind o metodă excelentă pentru investigarea modului în care sunt alcătuite și cum funcționează comunitățile);
- Biodiversitatea - Diversitatea ecologică;
- Analiza de asociere;
- Compararea sistemelor ecologice - Indici de similitudine;

- **Metode de ordonare și clasificare a sistemelor ecologice** - construirea dendrogramelor.

- Metode de studiu ale productivității și analiza rețelelor trofice... și multe altele...

Prin urmare tot ce conține acest capitol este pus într-un context clar; tema de față completează ceea ce (se presupune că) s-a învățat (cel puțin) la laboratoarele de ecologie și de modelare de până acum.

Comunitatea, definită ca un ansamblu de populații simpatrice și sincronice care interacționează pe mai multe niveluri, poate fi abordată teoretic și metodologic prin prisma a două puncte de vedere; ambele sunt valabile în anumite condiții și se completează reciproc:

- **Abordarea discontinuă** (a comunităților unitare, a cartării comunităților, abordarea integralistă *sensu* Clements) - comunitățile sunt identificabile, individualizate, delimitate, cu granițe bine definite între ele etc. Principalele metode de studiu au fost amintite anterior (cele de clasificare sunt în mod special folosite la această categorie) (*Discontinuum Approach*, în engleză) fiind totodată concepția clasică în fitocenotaxonomie (taxonomia asociațiilor de plante).

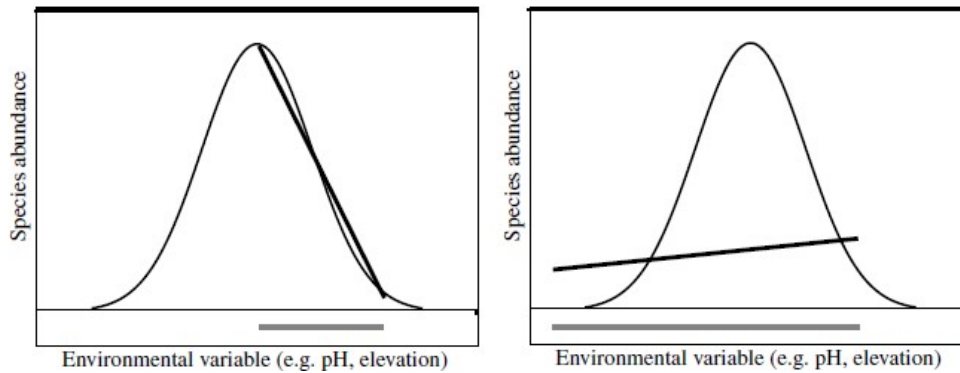
- **Abordarea continuă** (*Continuum Approach* în engleză), sau abordarea individualistă *sensu* Gleason, **analiza de gradient** - conform căreia nu există granițe distincte între comunități, se presupune că sunt zone largi de tranziție între diferitele sisteme ecologice adiacente etc. Metodele de ordonare prezentate aici sunt bazate pe ideea analizei de gradient a comunităților.

**O întrebare foarte importantă în analiza bazată pe abordarea continuă este:**

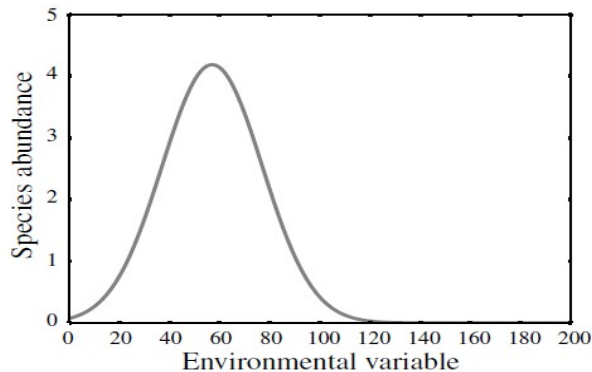
**Cum răspunde o specie (cum variază un parametru al acesteia) de-a lungul gradientului (variația treptată, continuă a) unui factor (variabilă, parametru) al mediului?**

- dacă segmentul de variație al gradientului este scurt, presupunem că abundența speciei variază liniar. Atunci aplicăm **modele liniare** (PCA, RDA);

- dacă segmentul este lung, presupunem că abundența speciei variază după o curbă unimodală (adică are o singură valoare maximă, un singur maxim sau mod), caz în care aplicăm **modele unimodale** (CA, CCA)



*Analiza de gradient pe baza modelului liniar, se pretează la segmente scurte ale domeniului de variație a variabilei de mediu (stânga sus), nu însă și dacă segmentul este lung (dreapta sus).*



*Răspunsul (și modelul) unimodal reprezentat mai sus, este mult mai corect atunci când domeniul de variație al factorului (a variabilei de mediu) delimitează sau include amplitudinea valențelor speciei (gradient lung).*

Răspunsul liniar al speciei pe un gradient scurt al variabilei de mediu este evaluat prin modelul liniar de regresie, respectiv prin metoda pătratelor minime. Pentru modelul unimodal de răspuns, cel mai simplu mod de a estima optimul speciei este de calcul a mediei ponderate ( $Mp(Sp)$ ) a valorilor variabilelor de mediu în cele  $n$  probe în care este prezentă specia. Valorile de abundență sau importanță a speciei sunt utilizate ca ponderi în calcularea mediei:

$$MP(Sp) = \frac{\sum_{i=1}^n X_i \cdot Abundenta_i}{\sum_{i=1}^n Abundenta_i}$$

unde  $n$  este numărul de probe,  $i$  este incrementul care ia valori succesive de la 1 la  $n$ ,  $X_i$  este valoarea variabilei de mediu  $X$  în proba  $i$ , iar  $Abundenta_i$  este valoarea abundenței speciei în proba  $i$ .

Când este necesar, toleranța speciei (lățimea curbei de tip clopot) poate fi calculată ca rădăcina pătrată a mediei ponderate a diferențelor pătrate dintre optimumul speciei și valoarea actuală dintr-o probă ( $SD(Sp)$ ). Valoarea este analogă cu abaterea (deviația) standard și servește ca bază pentru definirea unităților de SD pentru măsurarea lungimii axelor de ordonare.

$$SD(Sp) = \sqrt{\frac{\sum_{i=1}^n (X_i - MP(Sp))^2 \cdot Abundenta_i}{\sum_{i=1}^n Abundenta_i}}$$

În funcție de prezența sau absența variabilelor de mediu, precum și de câte variabile de specii sunt incluse în analize, metodele alese pot fi diferite.

### Tipuri de analize statistice

Variabila(e) de răspuns...(specii)	Variabila(e) de prognoză (de mediu)	
	Absentă	Prezente
... nu este niciuna	Atunci nu este nici ecologie și nici biologie -nu își are locul în acest capitol, volum și nici în disciplina studiată! Metoda: ați greșit facultatea și/sau specializarea!	
... este una	- rezumat de distribuție (al proprietăților de distribuție, adică grafice, frecvențe, histograme, statistică descriptivă etc.)	- modele de regresie <i>sensu lato</i>
... sunt mai multe	- analiza indirectă de gradient (PCA, CA, DCA, NMDS) - analiză de clasificare ierarhică (dendrograme)	- analiză directă de gradient (RDA, CCA) - analiza discriminantă (CVA)

Abrevierile provin din engleză:

PCA = "*principal component analysis*" (analiza componentelor principale)

DCA = "*detrended correspondence analysis*" (analiza de corespondență segmentată)

NMDS = "*non-metric multidimensional scaling*" (scalare multidimensională non-metrică)

RDA = "*redundancy analysis*" (analiza de redundanță)

CCA = "*canonical correspondence analysis*" (analiză de corespondență canonică)

CVA = "*canonical variate analysis*" (analiza variabilei canonice sau analiza funcției discriminante)



**Metodele de analiză a gradientului sau de ordonare** sunt de diverse categorii, cele mai importante fiind redate mai jos.

Nr. var. mediu	Nr. specii	Cunoștințe a priori asupra relației dintre specii și mediu?	Metoda	Rezultate
1, n	1	NU	Regresie	Dependența speciei de (relațiile cu) variabilele de mediu
Nici una	n	DA	Calibrare	Estimarea valorilor mediului
Nici una	n	NU	Ordonare necondiționată (independentă, ne-canonică)	Axele de variabilitate în compoziția speciilor; analiza indirectă de gradient
1, n	n	NU	Ordonare condiționată (constrânsă, dependentă, canonică)	Variabilitatea din compoziția speciilor explicată de variabilele de mediu. Relațiile variabilelor de mediu cu axele speciilor. Analiza directă de gradient.

**Tehnicile fundamentale de ordonare (analize și diagrame de ordonare)** se clasifică și se aleg în funcție de modul în care răspund speciile la variația gradientului (scurt sau lung, respectiv dacă sunt modele liniare sau unimodale) precum și dacă se consideră sau nu relația variabilelor dependente cu cele de mediu (canonice). Tabelul de mai jos reprezintă o sinteză a acestora.

	Metode liniare	Metode unimodale (medii ponderate)
<b>Necondiționate, non-canonice, neconstrânse (nu se ține seama sau nu avem variabile de mediu) = analiza indirectă de gradient</b>	Analiza componentelor principale ( <b>PCA</b> = <i>Principal components analysis</i> )	- Analiza de corespondență ( <b>CA</b> = <i>Correspondence analysis</i> ) - <b>DCA</b> = <i>Detrended correspondence analysis</i> )
<b>Condiționate, constrânse sau canonice (există și sunt considerate variabilele de mediu, independente) = analiza directă de gradient</b>	Analiza de redundanță ( <b>RDA</b> = <i>Redundancy analysis</i> )	Analiza de corespondență canonică ( <b>CCA</b> = <i>Canonical correspondence analysis</i> )

Ne oprim aici cu teoria, pentru a vedea cum se realizează o analiză de ordonare.

## Realizarea unei analize de ordonare, alcătuirea și interpretarea unei diagrame de ordonare în Canoco 4.5

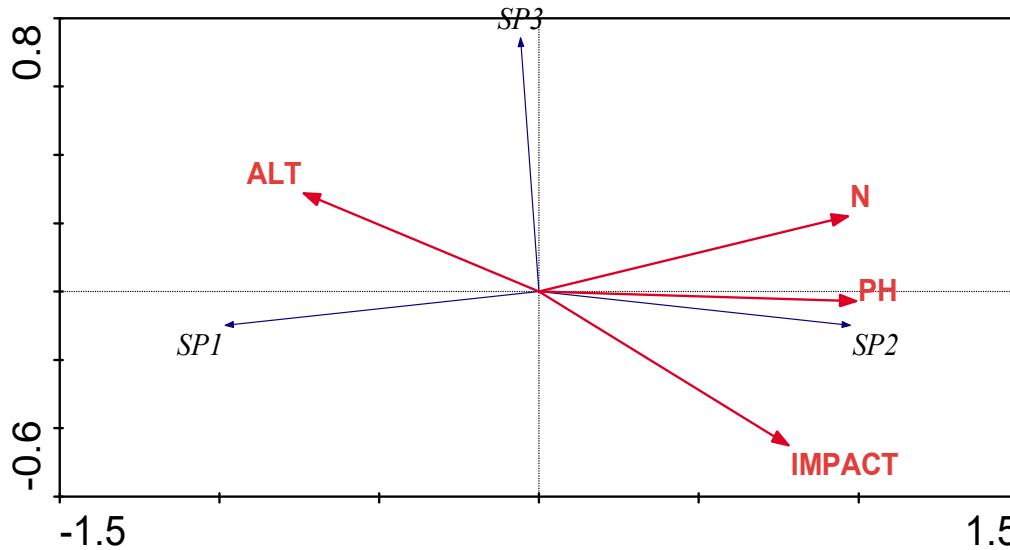
Un exemplu simplu (reduc ca număr de articole sau probe), prea puțin realist, dar sugestiv și un bun punct de plecare pentru a învăța organizarea datelor și analiza de ordonare a comunităților în *Canoco*.

Fie o comunitate de plante descrisă prin 9 probe (SAMPLE), 3 specii de cormofite (SP1, SP2, SP3) - acestea alcătuind variabilele care se referă la comunitate (variabile dependente sau "de specii"), iar separat scriem din nou probele (SAMPLE) și variabilele independente sau "de mediu": IMPACT - evaluarea mai mult sau mai puțin obiectivă a impactului antropic, pe o scară de la 1 (cel mai redus impact) la 6 (impactul cel mai puternic), PH - valoarea pH a solului din probe, N - conținutul de azot din sol și ALT - altitudinea la care se situează stația de eșantionare. Speciile sunt evaluate în termeni de abundență-dominanță relativă (sunt alese din asociații numai speciile cu abundență ridicată). Datele sunt introduse într-un program de foi de calcul (Excel sau altul):

	A	B	C	D	E	F	G	H	I	J
1	SAMPLE	SP1	SP2	SP3		SAMPLE	IMPACT	PH	N	ALT
2	1	35	2	51		1	1	2.1	12	563
3	2	23	4	32		2	3	2	13	612
4	3	21	3	24		3	2	2.5	15	620
5	4	34	5	56		4	3	2.1	11	510
6	5	12	6	21		5	4	3.2	23	462
7	6	2	45	16		6	5	6.7	56	374
8	7	0	56	58		7	4	7.9	89	451
9	8	4	45	35		8	6	5.3	45	304
10	9	11	4	27		9	3	2.8	22	450

Ce am dori să obținem în cele din urmă? De exemplu (ca mai jos) o diagramă care rezultă în urma unei analize de redundanță (RDA - *Redundancy Analysis*), în care toate datele sunt proiectate într-un spațiu 2D, definit de două axe care extrag cea mai mare parte din variabilitatea sau varianța datelor, săgețile albastre indicând sensul în care crește abundența speciilor, cele roșii indică sensul de creștere al valorilor variabilelor de mediu, iar direcțiile și unghiurile dintre săgeți, precum și dintre variabile și axe, indică gradul de corelație între diferitele entități.

Menționăm că această analiză și diagramă nu este obligatorie și nici singura posibilă!



Diagramă de ordonare (RDA) realizată pe baza datelor din fișierul Excel prezentat anterior.

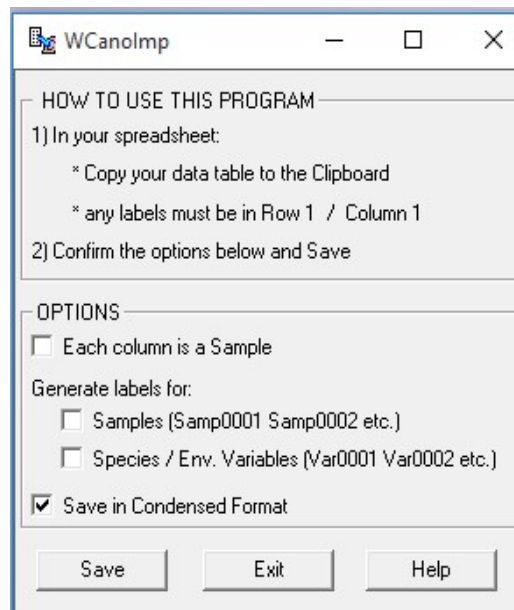
### **Rezolvare și (totodată) ghid succint de lucru în Canoco 4.5**

- Instalăm *Canoco* și așezăm (astfel încât să le avem la îndemână) icoanele pentru: introducerea datelor (*WCanoImp*), programul care realizează prelucrările și analizele (*Canoco for Windows*) și cel care realizează graficele (*CanoDraw for Windows*). Menținem deschis fișierul cu datele primare ale comunității și variabilele de mediu din programul respectiv (Excel sau altul).

Menționăm că această etapă este valabilă pentru versiunea 4.5 și cele inferioare. Începând cu *Canoco 5.x* toate aplicațiile sunt adunate (grupate) într-un singur program de bază, motiv pentru care se instalează numai acesta.



- Dublu clic pe *WCanoImp* → apare fereastra:



Faceți ce scrie acolo:

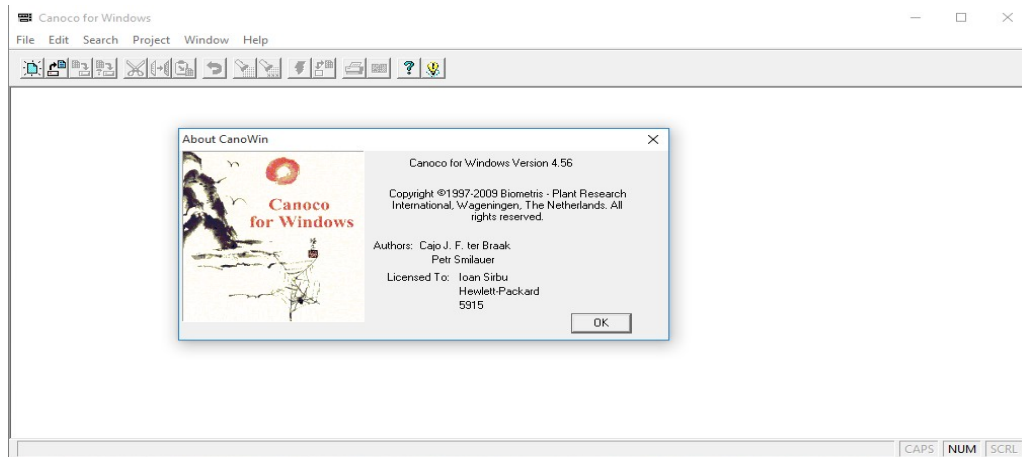
- Datele să fie astfel concepute încât pe prima linie să fie etichetele variabilelor, iar pe prima coloană să fie numărul curent al probelor sau al stațiilor de eșantionare.
- Dacă probele sunt pe coloane și speciile pe linii, selectați *Each column is a Sample*
- Copiați (Ctrl+C) sau puneți în Clipboard toate datele de specii (trasați un dreptunghi în jurul secțiunii fișierului care se referă la specii). Ignorați variabilele de mediu deocamdată.
- Clic pe Save.

	A	B	C	D	E	F	G	H	I	J
1	SAMPLE	SP1	SP2	SP3		SAMPLE	IMPACT	PH	N	ALT
2	1	35	2	51		1	1	2.1	12	563
3	2	23	4	32		2	3	2	13	612
4	3	21	3	24		3	2	2.5	15	620
5	4	34	5	56		4	3	2.1	11	510
6	5	12	6	21		5	4	3.2	23	462
7	6	2	45	16		6	5	6.7	56	374
8	7	0	56	58		7	4	7.9	89	451
9	8	4	45	35		8	6	5.3	45	304
10	9	11	4	27		9	3	2.8	22	450
11										

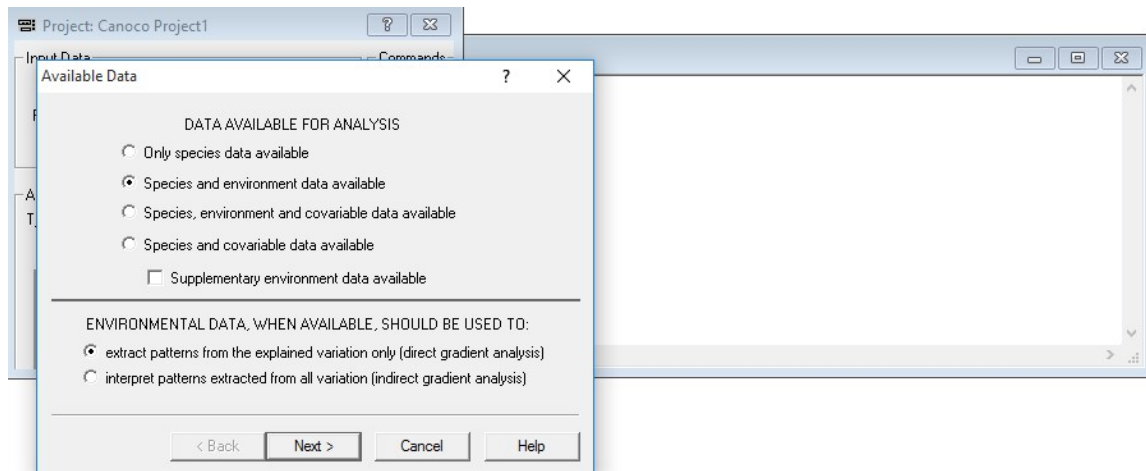
- Am marcat și am copiat (Ctrl+C)
- În *WCanolmp* alegem locația (directorul, fișierul) și dăm un nume sugestiv pentru acest fișier, cu extensia .dta. De exemplu: specii.dta. Scrierea extensiei nu este obligatorie!
- Clic pe Save
- Programul ne deschide o fereastră pentru a introduce, dacă dorim, un anumit titlu pe prima linie a acestui fișier (Title Line Text)

- Clic pe OK
- Programul ne anunță că a creat fișierul dorit.
- Apare fosta fereastră de introducere și salvare a fișierelor.
- Noi selectăm acum din fișierul de date primare variabilele și datele de mediu (SAMPLE, IMPACT, PH, B, ALT). Trasați un dreptunghi în jurul acestor date și variabile și copiați (Ctrl+C).
- Ne întoarcem la WCanImp în fereastra veche și clic din nou pe Save.
- Se deschide fereastra de salvare, scriem un nou nume cu extensia adecvată (de exemplu mediu.dta) și dăm clic pe Save. Suntem anunțați că fișierul a fost salvat.
- Clic pe OK și acum putem închide și fișierul de date și pe WCanImp (Exit).

**Deschidem *Canoco for Windows* (dublu clic pe icoană). Apare fereastra:**



- Selectăm: File, apoi New Project.
- Apare o fereastră care ne întreabă ce date sunt disponibile (Available Data)
- Într-un studiu real, mai întâi vom face câteva interogări numai pe datele de specii, pentru a afla informații despre acestea, precum și pentru a selecta tipul de analiză care se pretează la setul de date în cauză (vom afla dacă se poate realiza o analiză de redundanță RDA, sau una de corespondență canonică CCA, în funcție de lungimea segmentelor care rezultă dintr-o analiză preliminară DCA). Acum, pentru a simplifica, presupunem că am trecut de această etapă, care va fi învățată mai târziu, și realizăm o RDA. Vom testa setul de date și vom reprezenta grafic diagrama corespunzătoare.



- (acum, la început) selectăm ”Species and environment data available” (sunt disponibile date despre specii și mediu) din meniul de DATE DISPONIBILE PENTRU ANALIZĂ (DATA AVAILABLE FOR ANALYSIS)

- din meniul DATELE DE MEDIU, DACĂ SUNT DISPONIBILE, POT FI UTILIZATE PENTRU (ENVIRONMENTAL DATA, WHEN AVAILABLE, SHOULD BE USED TO), noi (în acest studiu de caz) vom alege să extragem modelul numai din variația explicată de către variabilele de mediu alese (adică analiza directă de gradient) - în original: *extract patterns from the explained variation only (direct gradient analysis)*.

- Cealaltă variantă (atunci când nu avem sau nu vrem să ținem seama decât de variabilele de specii, deci cele care se referă la comunitate), este de a selecta analiza indirectă de gradient (*indirect gradient analysis*).

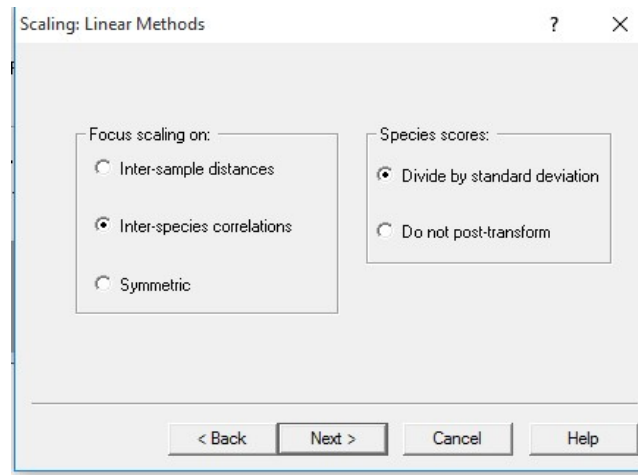
- Clic pe Next

- Programul ne solicită să încărcăm fișierele în următoarea fereastră. Cu Browse căutăm folderul unde avem datele și (sus) selectăm fișierul cu variabilele de specii (Species data file name), apoi (mai jos) fișierul cu date de mediu (Environment data file name), celelalte două câmpuri le lăsăm în pace, iar în ultima linie trebuie să alegem o cale și să scriem numele unui fișier care va conține soluțiile analizelor, fișier care trebuie să aibă extensia .sol. De exemplu vom alege numele exercitiu.sol și amplasăm noul fișier în aceeași locație cu fișierele de date.

- Clic pe Next

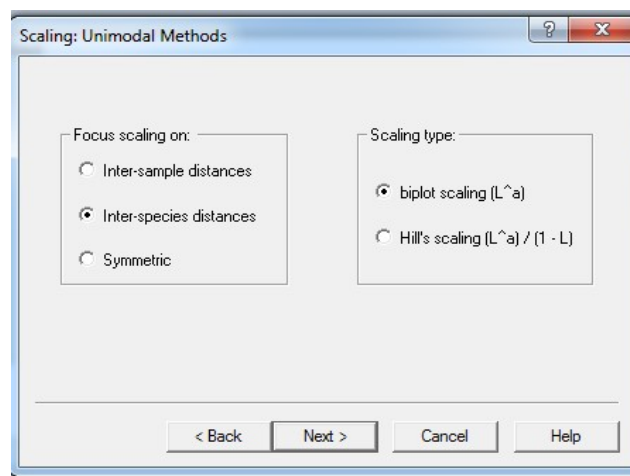
Response Models	Gradient Analysis Methods		
	Indirect	Direct	Hybrid
Linear	<input type="radio"/> PCA	<input checked="" type="radio"/> RDA	<input type="radio"/> hRDA
Unimodal	<input type="radio"/> CA	<input type="radio"/> CCA	<input type="radio"/> hCCA
Unimodal (detrended)	<input type="radio"/> DCA	<input type="radio"/> DCCA	<input type="radio"/> hDCCA

- Alegem tipul de analiză pe care dorim să o realizăm.
- Vom discuta mai târziu care sunt metodele, când și în ce condiții le alegem.
- Metoda de analiză a gradientului este indirectă, dacă nu avem decât variabile de specii (nu dispunem de date cu privire la mediu) sau nu vrem să le luăm în calcul, la început cel puțin. Analiza este directă în celălalt caz, dar există și posibilitatea să efectuăm o analiză hibridă.
- Deocamdată selectăm RDA (*Redundancy Analysis*) considerând că modelul de răspuns liniar este cel mai potrivit (ceea ce de fapt și este, vom vedea în continuare motivele).
- Next. Va aduce pe ecran **fereastra de scalare a scorurilor de ordonare**



- Deocamdată lăsați opțiunile din această fereastră așa cum sunt.

Ce sunt și cum utilizăm ferestrele de scalare? Mai sus am redat fereastra de scalare pentru metoda liniară. Pentru metodele unimodale fereastra de scalare este redată mai jos:



Precizia concluziilor cu privire la similaritatea probelor, a relațiilor dintre specii și/sau variabilele de mediu, depinde parțial și de scalarea relativă a scorurilor axelor individuale de ordonare. Opțiunile din ferestrele de mai sus permit alternative în interpretarea vizuală a rezultatelor ordonărilor, fără a afecta alte rezultate sau aspecte ale analizelor (cum ar fi valoarea variabilității explicate de fiecare axă, puterea și direcția relației între axe și variabile explicative, rezultatele testelor de semnificație sau de eroare de tip I în analizele canonice etc.). Când selectăm opțiunile trebuie să decidem dacă atenția se va concentra asupra probelor sau asupra speciilor în interpretarea diagramelor de ordonare. Apoi, la modele liniare, se ia decizia dacă lungimea săgeților trebuie să reflecte diferențele dintre abundențele speciilor (speciile mai abundente să aibă săgeți mai lungi) sau abundența speciilor individuale să fie transformată și raportată la o scară comparabilă. A doua opțiune corespunde la un *biplot de corelație*. Dacă (în modelul liniar) selectăm ca scorurile speciilor să fie

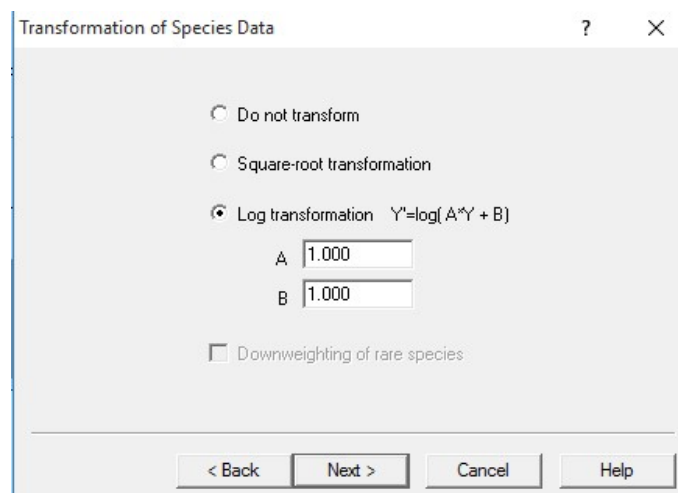


împărțite la abaterea standard (*Species scores: divide by standard deviation*), atunci lungimea fiecărei săgeți a unei specii exprimă cât de bine valorile speciei respective sunt approximate de către diagrama de ordonare. În cazul diagramelor de ordonare unimodală, vom alege scalarea Hill numai dacă gradientii compoziționali sunt foarte lungi (valoare mare a diversității beta dintre probe), în celelalte cazuri (și recomand utilizarea acestei opțiuni) scalarea de tip biplot oferă diagrame care pot fi interpretate într-o manieră cantitativă.

- Dați clic pe: Next.

- Teoria recomandă ca (în unele cazuri) datele comunității, care redau speciile în termeni de abundență relativă, să fie transformate. Se poate opta pentru transformare prin extragerea rădăcinii pătrate sau prin logaritmare. Dar, în cele mai multe date de comunități, sunt multe specii care apar rar (adică abundență zero în multe probe), motiv pentru care transformarea logaritmică nu se poate face direct (logaritm de zero este minus infinit). Variabilele se transformă după expresia:  $Y' = \log(A*Y + B)$ , unde - implicit - A și B sunt 1 (deci  $Y' = \log(Y + 1)$ ). Clic pe Next.

Motivul transformărilor este legat mai ales de relația dintre variabilele de mediu și cele de specii. În timp ce primele se modifică treptat, cele dependente prezintă o relație neliniară cu primele, de obicei. Multe valori ale speciilor în probe sunt zero, iar în majoritatea cazurilor valorile de abundență sunt fie mici, fie (atunci când speciile sunt în zona de optimum) foarte mari. Motivul pentru care datele de specii sunt transformate este de a defini o scară și o relație similară de variație cu cea a variabilelor de mediu, adică de a încerca identificarea unei relații de tip liniar (sau aproape) între cele două categorii. Transformarea logaritmică nu este singura posibilă: se poate opta și pentru transformarea prin extragerea radicalului sau alta.



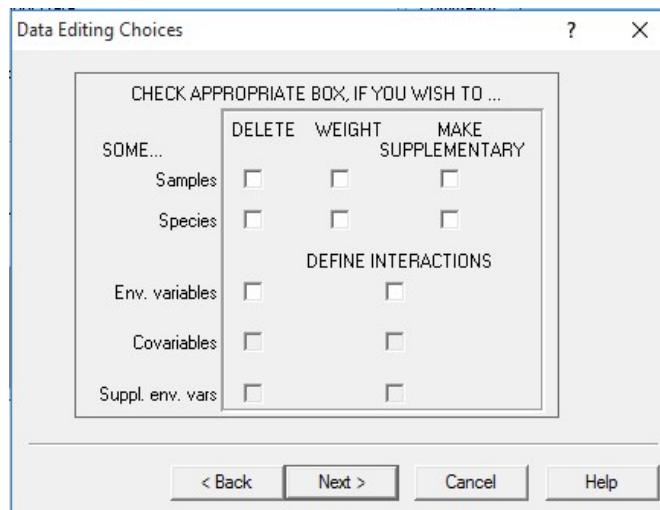
- Următoarea fereastră poate fi lăsată ca atare, și mergem mai departe (Next).

The image shows a dialog box titled "Centering and Standardization". It is divided into two main sections: "SAMPLES" on the left and "SPECIES" on the right. Each section contains four radio button options. In the "SAMPLES" section, the "None" option is selected. In the "SPECIES" section, the "Center by species" option is selected. At the bottom of the dialog, there are four buttons: "< Back", "Next >", "Cancel", and "Help".

Explicații suplimentare:

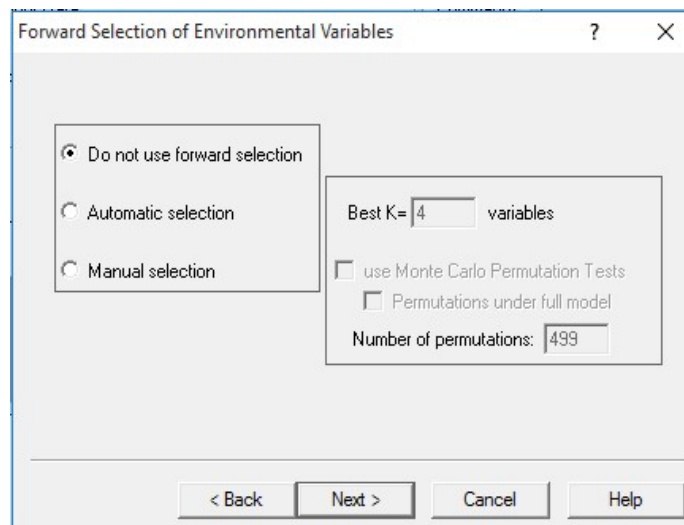
Această fereastră apare numai la metodele liniare (PCA sau RDA) și se referă la manipularea datelor din matricea comunității (variabilele de specii) înainte de calcularea ordonării. Centrarea pe probe (*Center by sample*) în coloana din stânga, imaginea de mai sus, produce o medie de zero pentru fiecare rând (articol). Centrarea pe specii (coloana din dreapta: *Center by species*) produce un rezultat al mediei de zero pentru fiecare coloană. Centrarea pe specii este obligatorie pentru metodele liniare canonice (RDA) sau pentru orice metodă liniară de ordonare parțială (adică acolo unde sunt utilizate covariate). Standardizarea pe specii (sau pe probe) produce ca norma fiecărei coloane (sau rând) să fie egală cu unu. Norma este rădăcina pătrată a sumei de pătrate ale valorilor din coloane (sau rânduri în cel de-al doilea caz). Dacă se optează atât pentru centrare cât și pentru standardizare, centrarea este realizată mai întâi. După centrarea și standardizarea pe specii, coloanele reprezintă variabile cu medie zero și varianță egală cu unitatea. Atunci, o PCA realizată pe datele de specii va corespunde la o "PCA pe o matrice de corelații" (între specii). Dacă nu se standardizează prin norma speciei, va rezulta o "PCA pe matrice de varianță-covarianță". De reținut faptul că standardizarea pe specii (și nu standardizare și nici centrare pe probe) trebuie realizată atunci când variabilele de specii (de răspuns) diferă în ceea ce privește scara de măsură. Dacă sunt variabile de mediu disponibile în metoda de ordonare (RDA), se poate selecta standardizarea prin eroarea varianței (*Standardize by error variance*), care va calcula, separat pentru fiecare specie, cât din varianța speciei nu a fost explicată de variabilele de mediu. Inversa erorii de varianță este **ponderea speciei**. Astfel, cu cât o specie este explicată mai bine de către variabilele de mediu disponibile, cu atât ponderea acesteia este mai mare în analiza finală.

- Next, va aduce fereastra opțiunilor de editare a datelor.



- Dacă avem opțiuni legate de editarea datelor, de a defini interacțiuni între variabile, de a stabili ponderi sau să eliminăm anumite articole și/sau variabile, putem să facem toate acestea și multe altele în această fereastră. Acum nu dorim să facem nimic de acest gen, așa că lăsăm în pace fereastra.

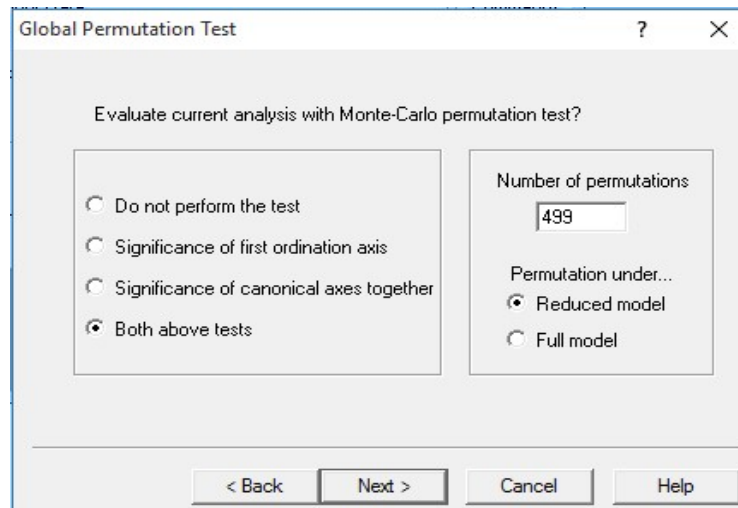
- Next.



- Același decizie. Implicit este selectată opțiunea ”Do not use forward selection”.

Dacă dorim să aflăm importanța fiecărei variabile de mediu, precum și a grupului de variabile, în ceea ce privește explicarea variabilelor dependente, vom selecta opțiunea *Automatic selection*, care va ordona variabilele de mediu în ceea ce privește importanța lor în explicarea variației compoziției speciilor.

- Next



- Mai importantă decât multe altele este testarea independenței dintre variabilele comunității (de specii) și cele de mediu, cu ajutorul testului de permutări prin metoda Monte-Carlo.

Din păcate am văzut că testarea aceasta este frecvent ignorată sau necunoscută, ceea ce nu este corect. De aceea putem alege ambele teste de semnificație (*Both above tests*), pentru a testa atât prima axă de ordonare cât și toate axele. În *Canoco 5* putem realiza câte un test separat pentru fiecare axă în parte, nu însă și în această versiune.

- Ipoteza nulă afirmă că nu există suficientă evidență pentru a susține o relație semnificativă (suficient de strânsă), de dependență, între variabilele de mediu și cele ale speciilor (altfel spus ***H<sub>0</sub>: variabilele independente (de mediu) nu explică variația compoziției comunităților***).

- Programul calculează un parametru statistic de tip F (Fischer) între cele două categorii de variabile, după care începe să facă permutări, adică amestecă seturile de date, la întâmplare. Aceeași linie (articol) de valori ale speciilor, este pusă în față, și comparată cu alte și alte seturi (linii, articole) de date ale variabilelor de mediu, care au caracterizat originar alte comunități, și se calculează în fiecare caz parametrul F al testului Fischer. Ideea de bază este că, dacă într-adevăr nu există relații între cele două categorii de variabile, atunci oricare set de date de specii poate fi explicat la fel de bine de oricare altă combinație de valori ale variabilelor de mediu.

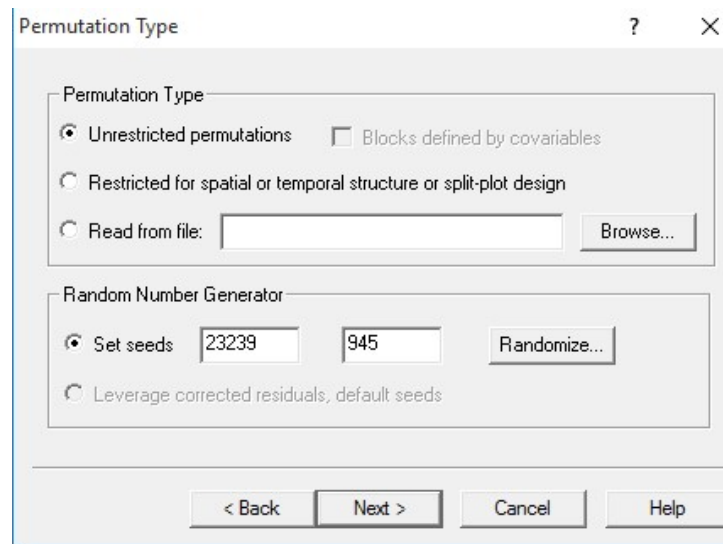
- Softul calculează de câte ori, în urma combinațiilor întâmplătoare, a rezultat o valoare a lui *F* mai mare decât cea calculată pe datele originale (setul inițial, real, de date).

$$p = \frac{n_x + 1}{N + 1}$$

unde:  $n_x$  = de câte ori (în urma câtor permutări) a rezultat o valoare  $F$  mai mare decât cea originală,  $N$  = numărul total de permutări,  $p$  = probabilitatea ca ipoteza nulă să fie adevărată.

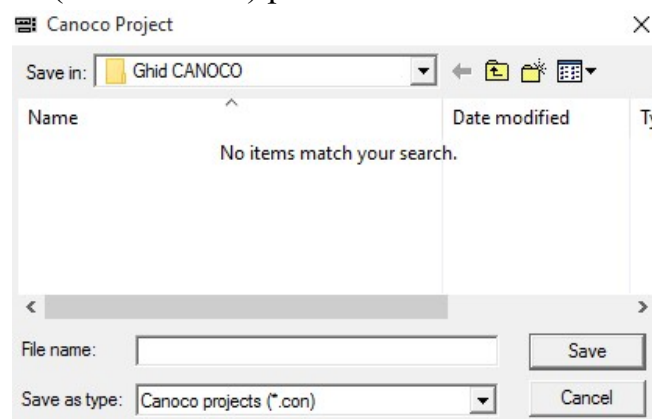
Dacă  $p < 0.05$  (sau nivelul de asigurare ales), respingem ipoteza nulă, caz în care afirmăm că există o relație statistic semnificativă între variabilele explicative și cele dependente (între cele de specii și de mediu), respectiv că variația compoziției comunităților este bine și semnificativ explicată de variabilele de mediu alese, în ceea ce privește (sau de către) o axă sau de toate axele de ordonare.

- Next.

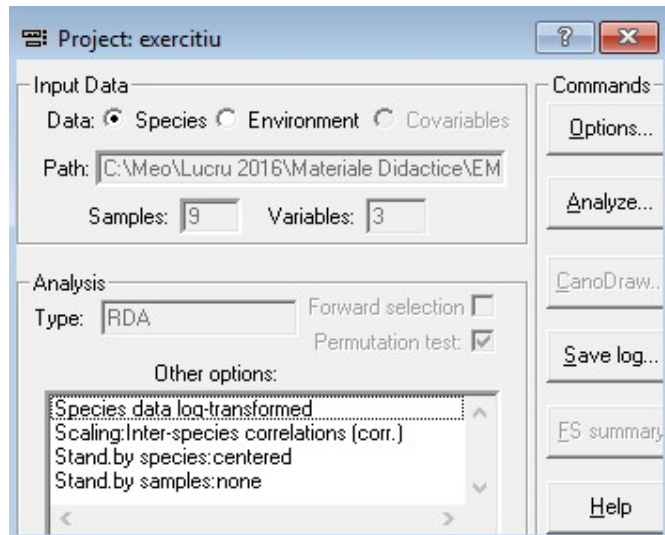


- Dacă dorim să modificăm sau să restricționăm permutările etc. putem să o facem din următoarea fereastră, dacă nu, lăsăm opțiunile implicite și trecem cu "Next" mai departe.

- Suntem anunțați că proiectul a fost încheiat, dar dacă mai dorim să revenim asupra unei opțiuni, acum este momentul, dacă nu, selectăm "Finish" și va apare fereastra de denumire a produsului (a rezultatului) proiectului nostru.



- Alegem un nume și eventual tastăm extensia pentru analiză sau proiect *.con*. De exemplu putem opta pentru *exercitiu.con*.
- Clic pe Save.



- Vedem din nou opțiunile și caracteristicile proiectului nostru, ne putem întoarce (Options) pentru a modifica proiectul, putem realiza analizele (Analyze), putem salva rezultatele (Save log) etc. Noi dăm clic pe Analyze și derulăm fereastra de rezultate (jurnalul).

```

**** Summary of Monte Carlo test ****

Test of significance of first canonical axis: eigenvalue = 0.897
                                           F-ratio   = 34.882
                                           P-value   = 0.0020

Test of significance of all canonical axes : Trace    = 0.948
                                           F-ratio   = 18.376
                                           P-value   = 0.0040

( 499 permutations under reduced model)

```

- Primim mai multe informații. La sfârșitul analizei (ca mai sus) aflăm dacă ipotezele nule care susțin independența dintre cele două categorii de variabile, respectiv faptul că prima și, separat, toate axele canonice, sunt semnificative. După cum vedem în ambele cazuri probabilitatea este foarte mică (sub 0.005), deci ipoteza nulă este respinsă. Acest fapt înseamnă că există relații semnificative între seturile noastre de variabile, deci studiul poate continua.
- Derulăm rezultatele analizei pentru a primi o sumedenie de informații, dar foarte importante, acum că știm că dispunem de un studiu (proiect) semnificativ, sunt cele redate astfel:

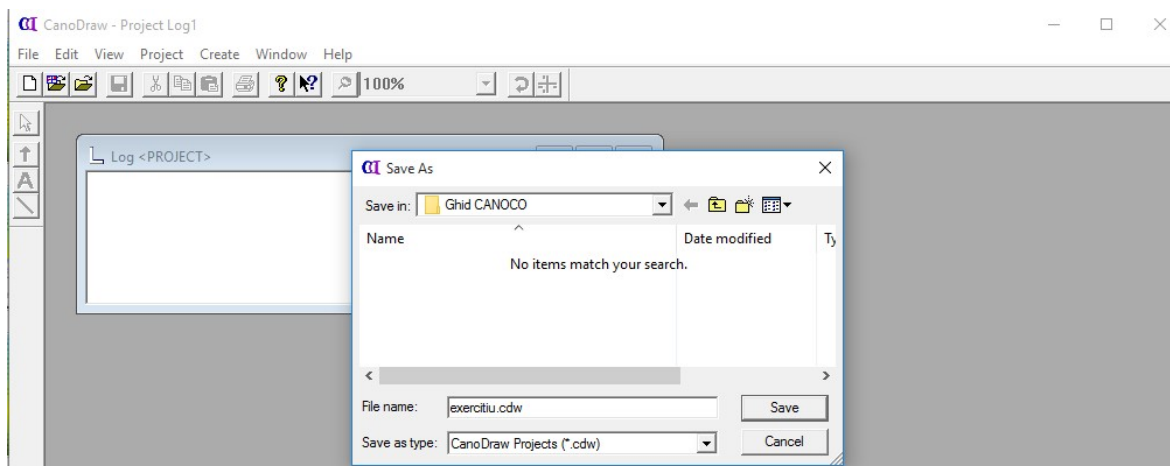
Axes	1	2	3	4	Total variance
Eigenvalues	0.897	0.047	0.004	0.050	1.000
Species-environment correlations	0.999	0.965	0.279	0.000	
Cumulative percentage variance					
of species data	89.7	94.4	94.8	99.8	
of species-environment relation:	94.6	99.6	100.0	0.0	
Sum of all eigenvalues					1.000
Sum of all canonical eigenvalues					0.948

Se observă din linia de valori ale proporției de varianță extrasă de către axe (respectiv din **valorile rădăcinilor latente** = "Eigenvalues" ale matricii) că prima axă este responsabilă pentru extragerea a 89.7% din variație, iar primele două axe explică 94.4% din toată variația conținută de către variabile. Prin urmare este foarte justificată reprezentarea variabilelor pe primele 2 axe.

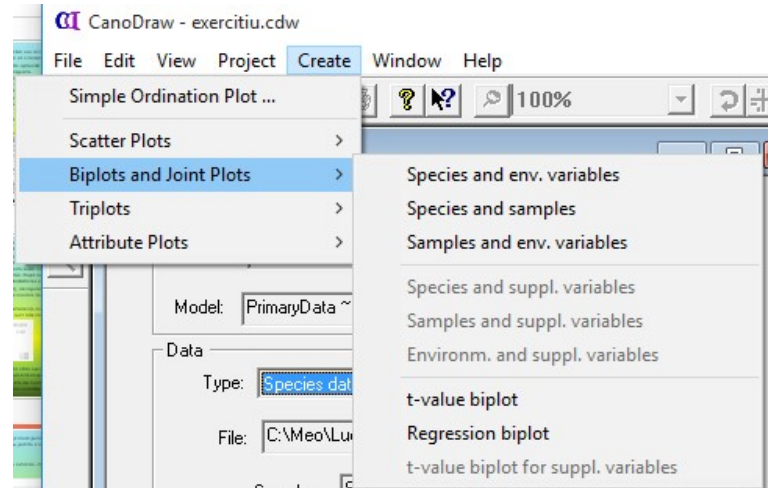
La întrebarea "câte procente din variația în compoziția comunității sunt explicate de variația variabilelor independente" (altfel spus care este echivalentul coeficientului de determinare din analiza de regresie pentru această analiză multivariată), răspunsul este dat de suma rădăcinilor latente canonice înmulțit cu 100 (*sum of all canonical eigenvalues* =  $0.948 * 100 = 94.8\%$ ). Prin urmare 94.8% din compoziția specifică a comunității este explicată de variabilele de mediu considerate. Concluzionăm că am obținut un model foarte bun!

- Ne întoarcem (*switch views*) la ultima fereastră a proiectului și acum putem selecta să intrăm în al treilea modul al acestui soft, și anume în *CanoDraw*, pentru a reprezenta relațiile și diagramele dorite. Clic pe *CanoDraw*.

- Vom alege un nou nume pentru fișierul din modulul grafic, cu extensia .cdw. De exemplu *exercitiu.cdw* (ca în fereastra de mai jos). Clic pe *Save*.



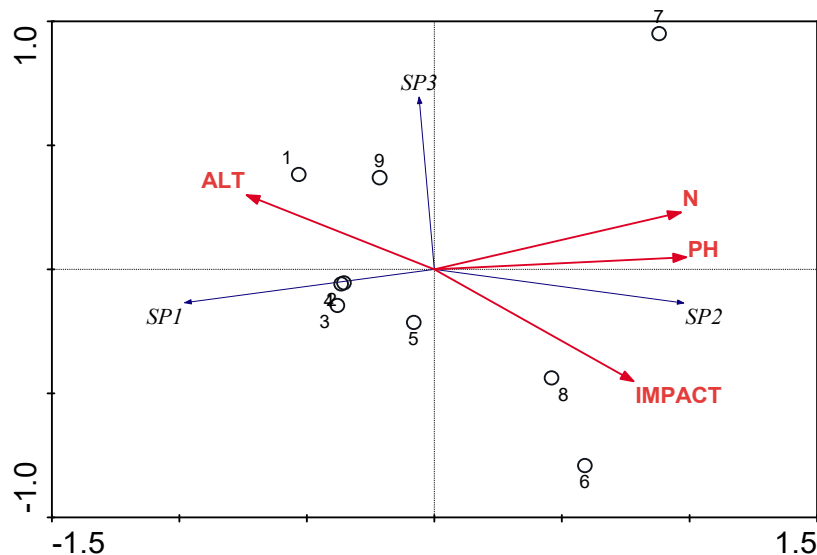
- Și acum (în sfârșit) putem realiza diagrame de ordonare și multe alte grafice. De asemenea putem selecta diverse opțiuni și condiții de reprezentare grafică (putem include sau exclude variabile și specii, de exemplu cele foarte rare etc.)



- Prin selecția: Create → Biplots and Joint Plots → Species and env. variables - rezultă diagrama de ordonare care corespunde unei analize de redundanță (RDA), adică cea pe care am dorit să o realizăm;

- urmează interpretarea...

- sau putem să selectăm Triplots și apoi clic pe "with environmental variables", obținând:



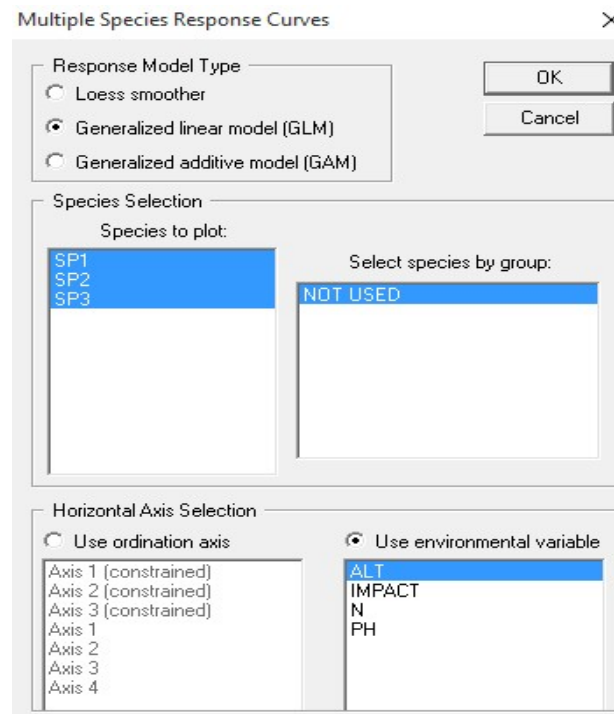
În aceasta vedem poziția relativă a speciilor, a variabilelor de mediu (independente sau explicative), dar și probele (stațiile de eșantionare) sau comunitățile



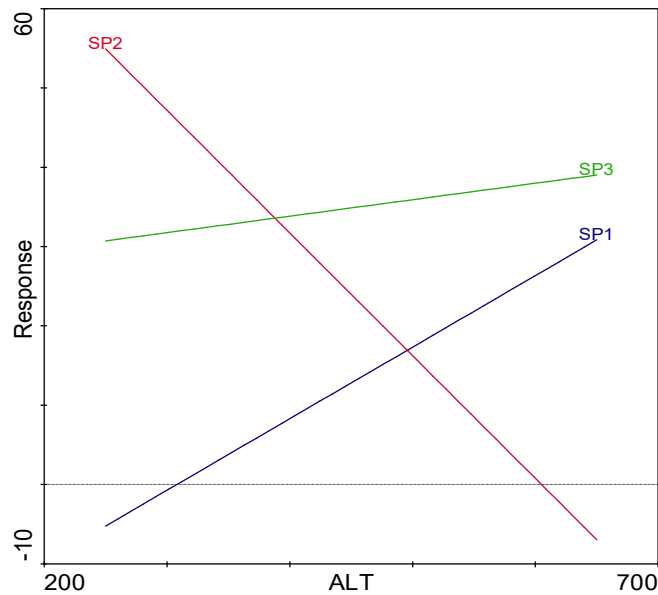
individuale și caracteristicile mediului lor. Această analiză permite interpretări de amănunt cu privire la preferințele speciilor pentru diferitele valori (domenii de variație) ale variabilelor de mediu, asemănări și deosebiri între valențele ecologice dar și între dimensiunea de habitat a nișelor ecologice etc.

Răspunsul speciilor la o anumită variabilă se face prin analiza de regresie; *Canoco* permite o serie de opțiuni pentru aceasta (regresie liniară, pătratică, cubică etc.), precum și testarea semnificației acestora. De exemplu dacă dorim să vedem răspunsul după un model liniar al variației abundenței speciilor la altitudine, selectăm: Create → Attribute Plots → Species response curves...

urmează alte câteva opțiuni ....



... și (după câteva ferestre de rezultate) rezultă graficul de mai jos. Din rezultatele intermediare aflăm că relațiile dintre abundența speciei 1, a cărei abundență crește cu altitudinea și cea a speciei 2, care descrește cu altitudinea, sunt semnificative, nu însă și a speciei 3 care pare să fie independentă de aceasta.



**Cum selectăm metoda de analiză de ordonare (model liniar sau unimodal)?**

**Metodele liniare sunt:**

- potrivite pentru analiza gradientilor scurți de mediu,
- pentru analiza experimentelor manipulative cu mulți factori,
- analiza seturilor de date omogene.

**Metodele unimodale sunt:**

- potrivite pentru analiza gradientilor lungi de mediu,
- a seturilor de date heterogene,
- a variabilelor numerice continue de mediu.

**Algoritmul simplificat pentru selecția metodei:**

1. Rulăm inițial o DCA (*detrended correspondence analysis*) și ne uităm la lungimea cea mai mare a gradientilor (căutăm și notăm cea mai mare valoare a *Length of gradients*);
2. Dacă lungimea maximă este **sub valoarea 3, utilizați metode liniare**;
3. Dacă lungimea este **peste 4, utilizați metode unimodale**;
4. Dacă lungimea este **între 3 și 4, puteți utiliza oricare dintre metode**.

Aplicația DCA și evaluarea lungimii gradientului primei axe de ordonare la exercițiul nostru; observăm că lungimea este de 0.858, mult sub valoarea de 3, motiv pentru care am ales o analiză RDA și nu CCA!

Type of Analysis ? X

Gradient Analysis Methods

Response Models	Indirect	Direct	Hybrid
Linear	<input type="radio"/> PCA	<input type="radio"/> RDA	<input type="radio"/> hRDA
Unimodal	<input type="radio"/> CA	<input type="radio"/> CCA	<input type="radio"/> hCCA
Unimodal (detrended)	<input checked="" type="radio"/> DCA	<input type="radio"/> DCCA	<input type="radio"/> hDCCA

< Back Next > Cancel Help

---

Detrending Method ? X

Select method of detrending:

by segments

by 2nd order polynomials

by 3rd order polynomials

by 4th order polynomials

< Back Next > Cancel Help

```

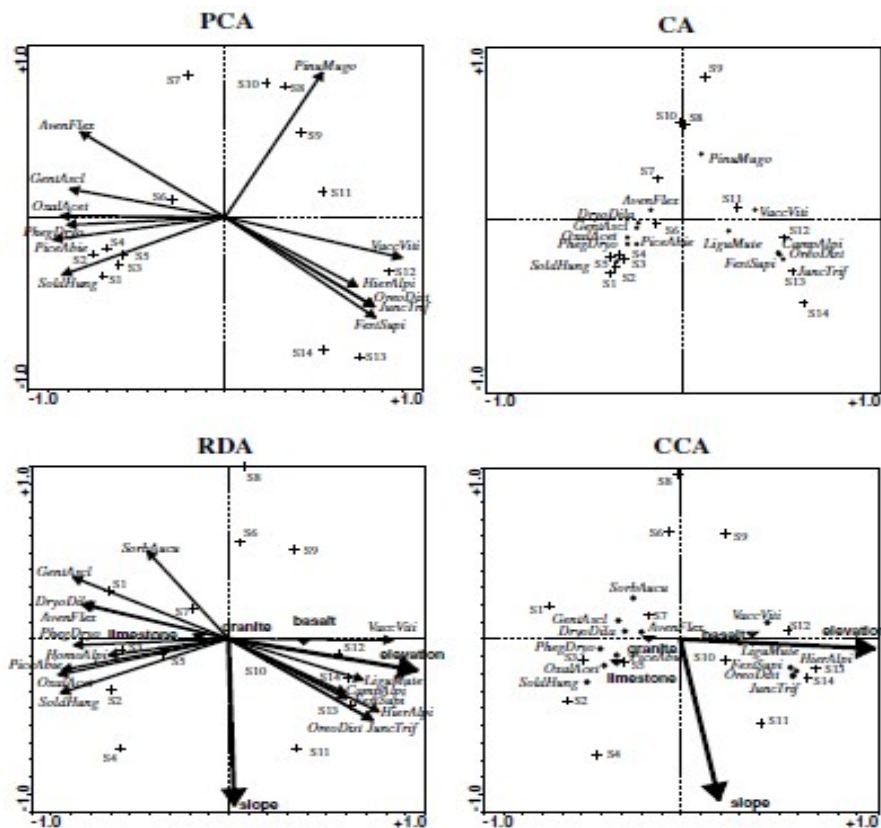
Log: exercitiu
**** Summary ****

Axes                1      2      3      4      Total inertia
Eigenvalues          : 0.127 0.001 0.000 0.000      0.133
Lengths of gradient : 0.858 0.610 0.000 0.000
Cumulative percentage variance
  of species data   : 95.9 96.3 0.0 0.0

Sum of all          eigenvalues      0.133
[Wed Jan 27 00:55:27 2016] CANOCO call succeeded

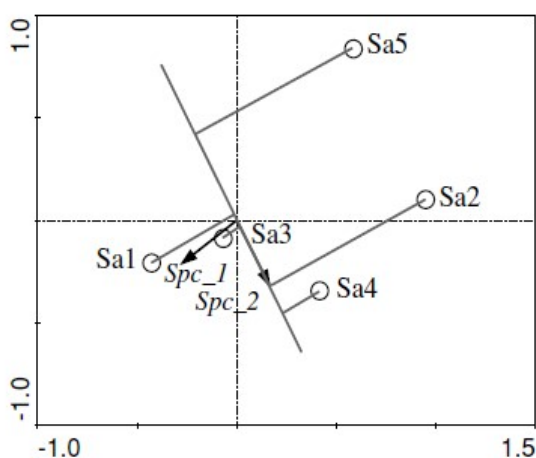
```

## Exemple de diagrame de ordonare



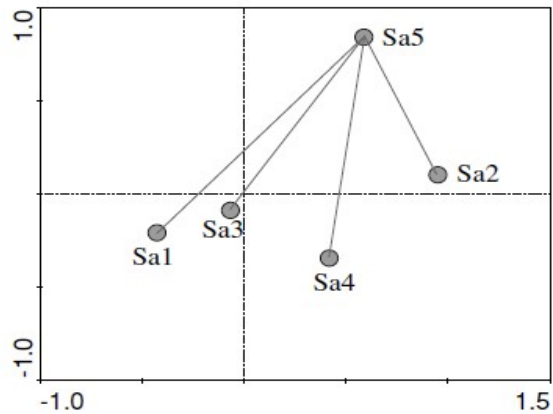
### Diagrame de ordonare și interpretarea acestora - metode liniare

Mai jos redăm diverse scheme de diagrame de ordonare, sub care specificăm câteva reguli de interpretare.

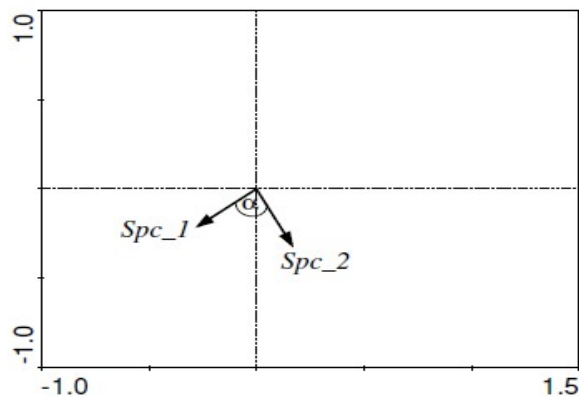


Proiectarea probelor pe vectorii speciilor (pe *Spc\_2*) într-un biplot dintr-o diagramă de ordonare liniară. Aici prognozăm cea mai mare abundență a *Spc\_2* în probele Sa4

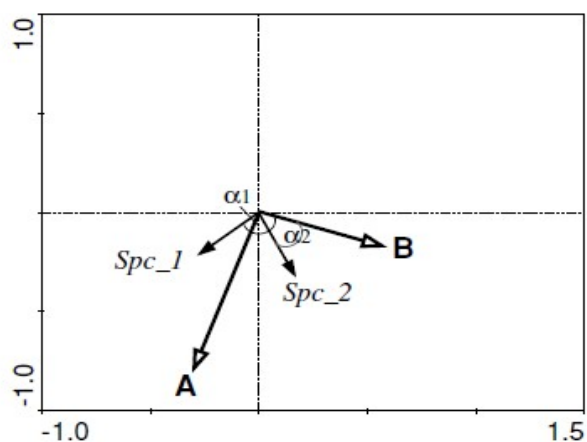
urmată de Sa2, apoi în Sa3 abundența speciei este aproape de valoarea medie din probe (proiecția este în zona de origine), valoarea cea mai mică fiind în Sa5.



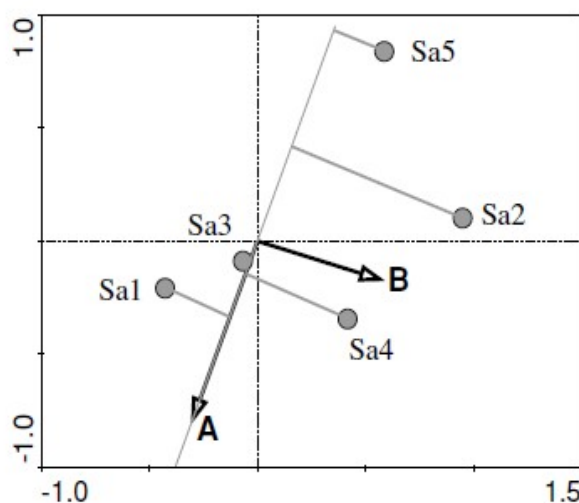
Distanța dintre probe într-o diagramă de ordonare. Dacă măsurăm disimilaritățile dintre proba Sa5 și celelalte, utilizând distanța euclidiană, cea mai scurtă distanță este între Sa5 și Sa2 (sunt cele mai asemănătoare) iar cea mai mare distanță este între aceasta și Sa1, deci prognozăm că sunt cele mai disimilare (deci diferite ca și compoziție a speciilor).



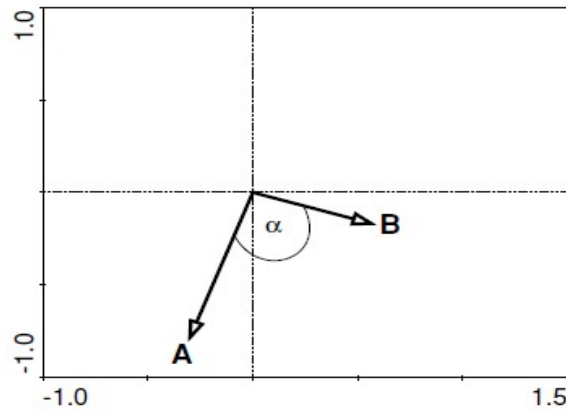
Unghiuri dintre specii într-o diagramă realizată într-o analiză de ordonare liniară. Deoarece săgețile celor două specii sunt la un unghi de aproape  $90^\circ$ , între cele două specii se prognozează că nu există nici o corelație (aproape zero), adică sunt independente în spațiul de ordonare.



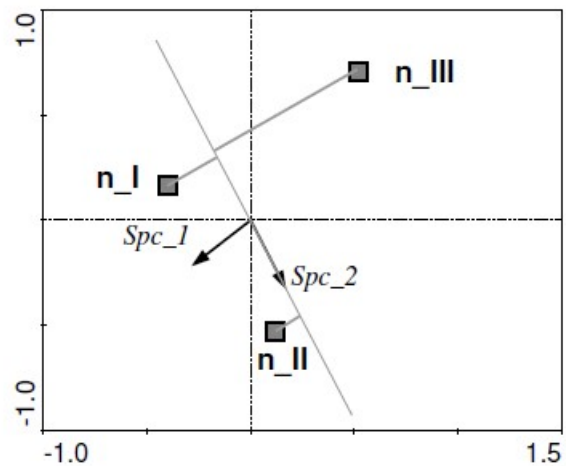
Unghiuri între specii și variabilele de mediu într-o diagramă de ordonare prin model liniar. Unghiul speciei *Spc\_1* este peste  $90^0$  fapt care indică o relație negativă (abundența speciei 1 scade cu creșterea lui B, deci regresie negativă), în timp ce unghiul mic (ascuțit) cu B al lui *Spc\_2* indică faptul că *Spc\_2* tinde să crească în abundență atunci când B crește. **Regula biplotului (proiectarea vârfulor săgeților speciilor pe săgețile variabilelor de mediu)** oferă o aproximare și mai precisă.



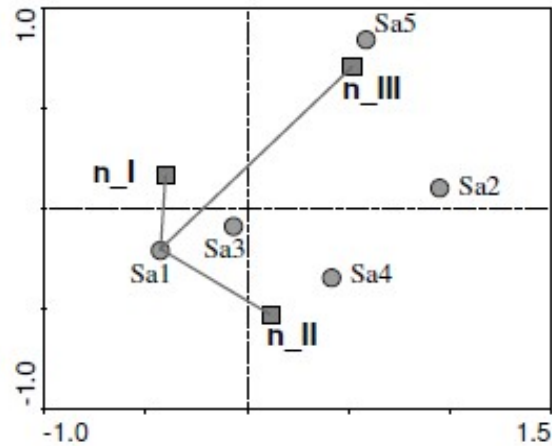
Proiectarea probelor pe săgețile variabilelor de mediu cantitative. Observați că variabila A este prognozată să prezinte valori similare pentru probele Sa3 și Sa4 deși acestea sunt localizate la distanțe și direcții diferite de linia lui A. A prezintă cele mai mari valori în proba 1 și cele mai mici în proba 5.



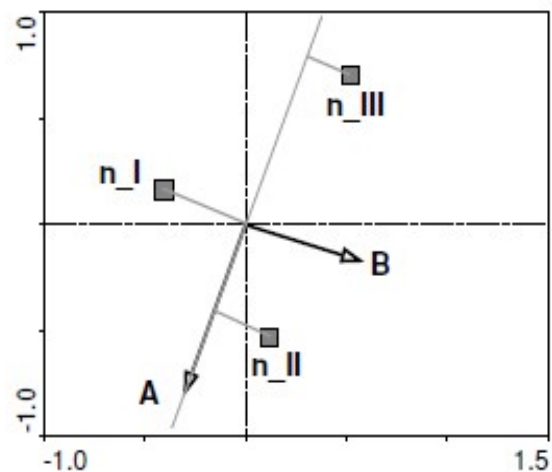
Măsurarea unghiurilor dintre două variabile de mediu. Unghiul aproape drept sugerează că între cele două variabile nu există aproape nici o corelație (valoare aproape de zero).



Proiectarea centroizilor variabilelor substitutive (*dummy variables*) pe săgeți ale speciilor într-o diagramă provenită dintr-o metodă liniară. Prognozăm că cea mai mare abundență medie a speciei *Spc\_2* este în probele care aparțin clasei *n\_II*, în timp ce în cele din clasele *n\_I* și *n\_III*, specia respectivă are o abundență medie prognozată mai mică și similară.



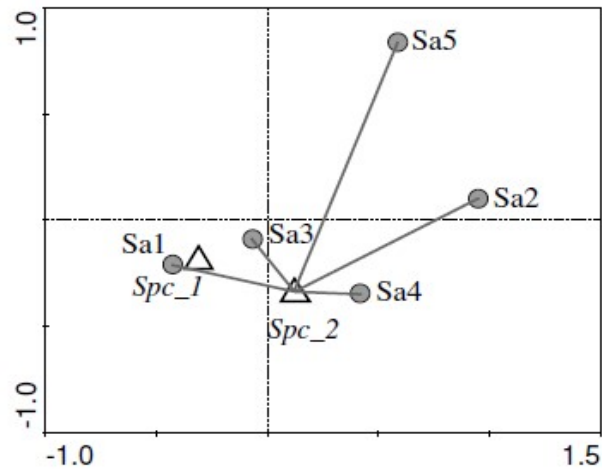
Măsurarea distanței dintre probe și centroizii variabilelor substitutive. Aici prognozăm că proba Sa1 are cea mai mare probabilitate să aparțină clasei n\_I și cea mai mică de a aparține clasei n\_III. Aceleași interpretări le vom realiza și în versiunile superioare ale *Canoco 5.x* unde se lucrează cu niveluri ale variabilelor factoriale și nu cu variabile substitutive.



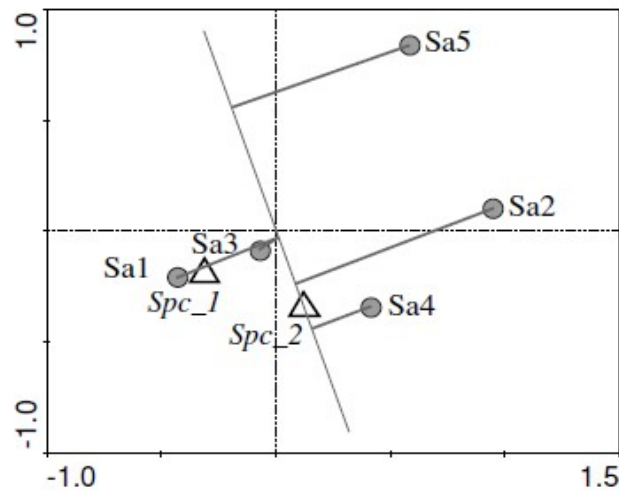
Proiectarea centroizilor variabilelor substitutive pe săgețile variabilelor cantitative de mediu. Probele aparținând clasei n\_II sunt prognozate să aibă cea mai mare valoare a variabilei A, urmate de cele din clasa n\_I și apoi de cele din n\_III în care A are cea mai mică valoare medie prognozată.



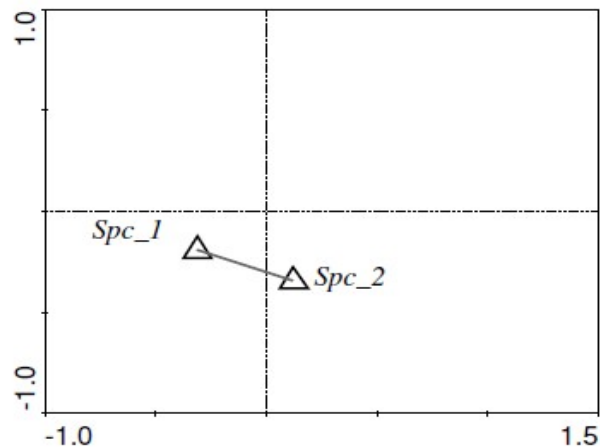
## Metode și diagrame unimodale de ordonare



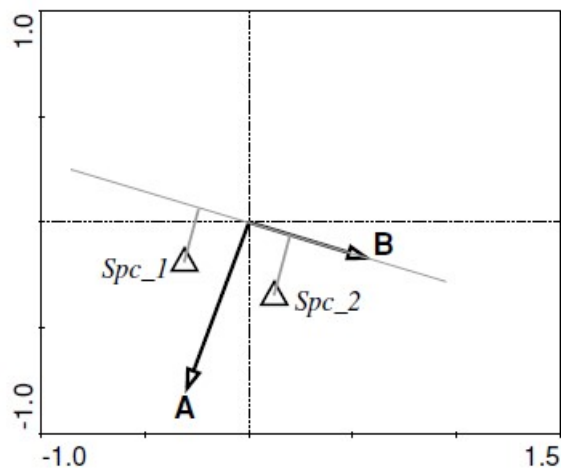
Distanța dintre specii și probe într-o diagramă unimodală. Specia Spc\_2 are probabil cea mai mare abundență relativă în probele Sa4 și Sa3 și cea mai mică abundență în Sa5 (probabil lipsește din aceasta).



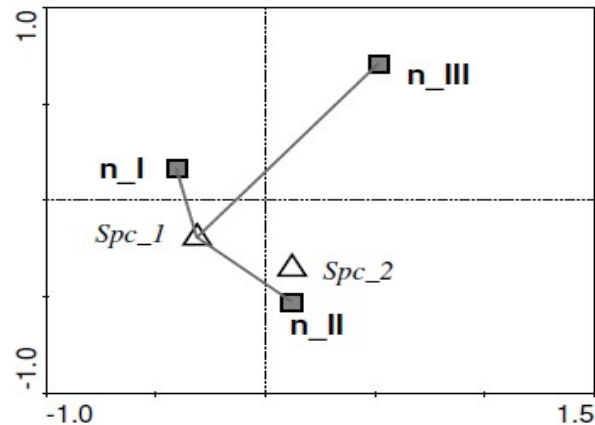
Regula biplotului aplicată la specii și probe în diagrama de ordonare provenită dintr-un model unimodal (o metodă unimodală de ordonare). Specia Spc\_2 are prognozabil cea mai mare frecvență relativă în probele Sa4 și Sa2 și cea mai mică în Sa5. Această specie este prognozată să apară în probele Sa1 și Sa3 la valoarea medie a frecvenței sau a abundenței ei relative. Pentru fiecare specie se poate trasa o linie (imaginară) între punctul de optim și origine, în scopul facilitării aplicării regulii de proiecție a biplotului.



Măsurarea distanțelor dintre specii într-o diagramă de ordonare realizată printr-o metodă unimodală. Distanța dintre specii într-o scalare biplot (concentrată pe distanța dintre specii) aproximează distanța chi-pătrat dintre distribuția speciilor.



Proiecția speciilor pe săgețile variabilelor cantitative de mediu. Interpretăm graficul de mai sus afirmând că specia Spc\_2 are probabil optimumul față de variabila B la valori mai mari ale acesteia decât specia Spc\_1. Proiectând speciile pe săgețile variabilelor de mediu putem să le ordonăm în funcție de valoarea optimă a acestora pe gradientul variabilei. În mod similar putem proiecta probe pe săgețile variabilelor explicative pentru a aproxima valorile din tabelul datelor de mediu.



Măsurarea distanțelor dintre specii și centroizii variabilelor substitutive de mediu. Frecvența medie relativă a speciei *Spc\_1* este prognozată de a fi maximă în probele din clasa *n\_I*, urmată de cele din clasa *n\_II* și minimă în clasa *n\_III*.

### Studiu de caz 2: comunități de păsări dintr-o pădure

Scop: studiul variabilității comunităților de păsări în relație cu variația variabilelor de habitat (variabilitatea mediului)

Metodă: s-au numărat perechile de păsări cuibăritoare, prin ascultarea cântecelor specifice, de 4 ori, câte două investigații (campanii de teren) executate în două sezoane de reproducere succesive, în cadrul unei rețele de puncte echidistante, alese într-o pădure montană. Unele puncte au fost amplasate într-un arboret dominat de molid, altele de fag. Variația mediului se referă la gradul de acoperire cu vegetație (lemnoasă, tufișuri, pășuni instalate în urma defrișării), altitudine, pantă, expoziție etc. (studiu realizat de M.E. Šálek în Munții Velka Fatra, Slovacia).

Variabile: încărcați fișierul *birds.xls* care conține două tabele (fișiere) - **fișierul se găsește pe pagina materiei.**

- Variabilele dependente (specii) sunt reprezentate prin valorile medii ale abundenței (37 de specii respectiv variabile), respectiv suma perechilor (cântecelor) numărate în fiecare pătrat, împărțit la 4 (cele patru studii de teren);
- Variabilele de mediu (care descriu caracteristicile habitatelor) sunt amplasate separat, dar în același document (foaie de calcul) și sunt 13 la număr.

#### Variabilele de mediu sunt:

- *Altit* = singura valoare numerică, reală, semnificând media altitudinii pătratului de eșantionare (sunt 43 de pătrate corespunzătoare);
- Celelalte entități conțin un amestec de variabile strict ordinale sau semi-cantitative;
- *Forest* = acoperirea suprafeței pătratului de eșantionare cu arbori (megafanerofite);
- *ForDens* = densitatea medie a vegetației arborescente;

- **BrLeaf** = frecvența relativă a foioaselor în arboret (0 = numai molidiș, până la 4 = numai foioase în stratul de arbori)
- **E2** = acoperirea cu tufărișuri;
- **E2Con** = procentul de conifere (molid) în stratul de tufărișuri;
- **E1** = acoperirea cu stratul ierbos și **E1Height** = înălțimea medie a acestuia;
- **Slope** = panta în grade împărțite la 5;

Urmează două perechi de variabile substitutive (*dummy variables*) care codifică două niveluri posibile pentru doi factori: poziția unor stânci mari pe suprafața pătratului **Rocks** și **NoRocks** precum și poziția sau nu a pătratului pe pante însoțite, adică sud-est, sud și sud-vest = **Warm**, respectiv **Cold** în celelalte cazuri.

### Mod de lucru:

- Introducem fișierele (variabilele și valorile corespunzătoare) în WCanoImp (putem în format întreg sau condensat) prin marcarea și copiere, alegem numele lor, ex.:

**birds\_spe.dta** pentru datele de specii și

**birds\_env.dta** pentru datele de mediu.

- Putem selecta *Species and environmental data available*, dar mai întâi vom decide asupra tipului de model pe care îl vom aplica, motiv pentru care în prima fereastră alegem:

- *indirect gradient analysis*. În fereastra următoare selectați fișierele și generați numele fișierului de soluții (de exemplu **dca.sol**)

În ferestrele următoare vom selecta:

DCA > *by segments* în pagina de *Detrending Method* > *Log transformation* în următoarea > lăsați în pace celelalte ferestre clic pe *Finish* și dați un nume proiectului (ex. **dca.con**) și

cereți analiza (*Analyze*)

- treceți în LogView, observați că primele două axe explică cca. 36% din variabilitatea în datele de specii, iar în linia care redă lungimea gradientilor, vedem cifrele:

*Lengths of gradient* : 2.001 1.634 1.214 1.613

- Gradientii sunt relativ scurți, motiv pentru care optăm pentru un model liniar.

- Începem un nou proiect pe baza celui precedent, deoarece o serie de opțiuni rămân valabile.

- Cu proiectul **dca.con** deschis în spațiul de lucru, selectăm *File* > *Save As...* din meniul de comenzi, după care scriem un nou nume (de exemplu **rda.con**). Când solicită programul, vom șterge informațiile din fereastra de rezultate (jurnalul); clic pe *Clear the log window*.

- Clic pe *Options* apoi lăsăm *Species and environment data available* dar apoi selectăm: *direct gradient analysis*

- Vom introduce, atunci când softul o cere, un nou nume pentru fișierul de soluții

(de exemplu *rda.sol*).

- În fereastra de *Type of Analysis* selectăm **RDA**;

În rest, lăsăm neschimbate toate celelalte selecții implicite, deoarece acestea sunt potrivite acestei analize.

- În fereastra de *Global Permutation Test* selectăm *Both above tests*;

- Solicităm analiza și ne uităm la valorile afișate în fișierul de rezultate;

- Testele indică respingerea ipotezelor nule ( $p=0.002$ ), deci există relații semnificative între variabilele dependente și cele de mediu (putem merge mai departe);

\*\*\*\* Summary \*\*\*\*

Axes		1	2	3	4	Total variance
Eigenvalues	:	0.170	0.088	0.053	0.041	1.000
Species-environment correlations	:	0.919	0.846	0.794	0.821	
Cumulative percentage variance						
of species data	:	17.0	25.8	31.0	35.2	
of species-environment relation:		38.4	58.2	70.1	79.5	
Sum of all eigenvalues						1.000
Sum of all canonical eigenvalues						0.442

All four eigenvalues reported above are canonical and correspond to axes that are constrained by the environmental variables.

Constatăm că variabilele de mediu alese explică 44.2% (sum of all eigenvalues \* 100) din variația compoziției în specii.

- fără a mai intra în alte amănunte, deschidem CanoDraw, alegem un nou nume pentru grafic (de exemplu *rda.cdw* și trebuie să declarăm variabilele de tip indicator sau de substituție);

- specificăm variabilele nominale *Rocks*, *NoRocks*, *Warm* și *Cold* (utilizând calea comenzilor: *Project > Nominal variables > Environmental variables*);

- solicităm graficul pe calea:

*Create > Biplots and Joint Plots > Species and env. variables*

Rezultatul o să ni se pară foarte încărcat:

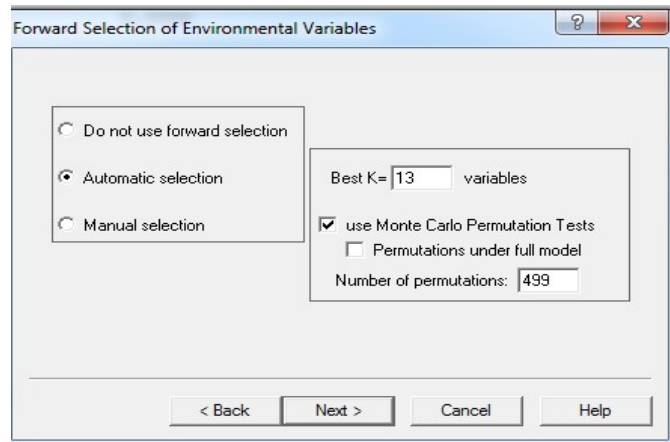


Rezultă o nouă diagramă RDA. Graficul este mai "aerisit": cum interpretați? (activitate independentă)

- **Cum clasificăm variabilele de mediu în ceea ce privește contribuția lor la rezultate, respectiv la explicarea compoziției specifice a comunităților de păsări?**

Rezolvare:

Solicităm o nouă analiză (Options), și trecem de ferestrele de opțiuni până la fereastra *Forward selection of environmental variables*, unde selectăm opțiunea *Automatic selection*.



Dăm clic pe Next și apoi selectăm opțiunea *Unrestricted permutation* în fereastra de permutări. Urmează clic pe Finish apoi pe Analyze și vom selecta din ultima fereastră opțiunea "FS summary" care va aduce un nou ecran, ca mai jos:

Marginal Effects				
Variable	Var.N	lam...		
Altit	1	0.11		
BrLeaf	4	0.10		
Forest	2	0.09		
E2Con	6	0.07		

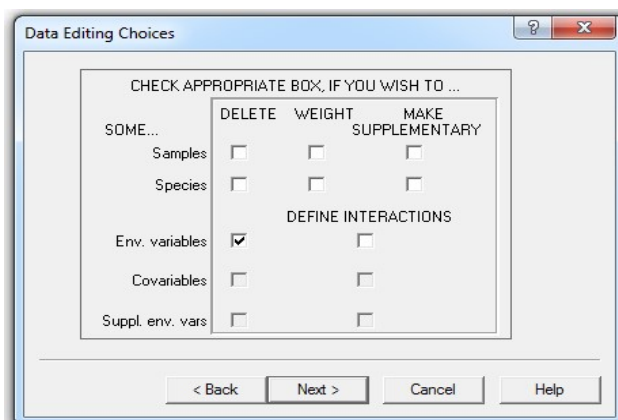
Conditional effects				
Variable	Var.N	lam...	P	F
Altit	1	0.11	0.002	5.32
BrLeaf	4	0.08	0.002	3.72
Forest	2	0.06	0.002	2.91
Warm	12	0.03	0.032	1.65

Vom da clic pe *Copy* și vom vizualiza (Ctrl+V sau Paste) rezultatele într-un document Word sau într-un fișier de foi de calcul, sau în Excel. Vom obțien tabelul de mai jos:

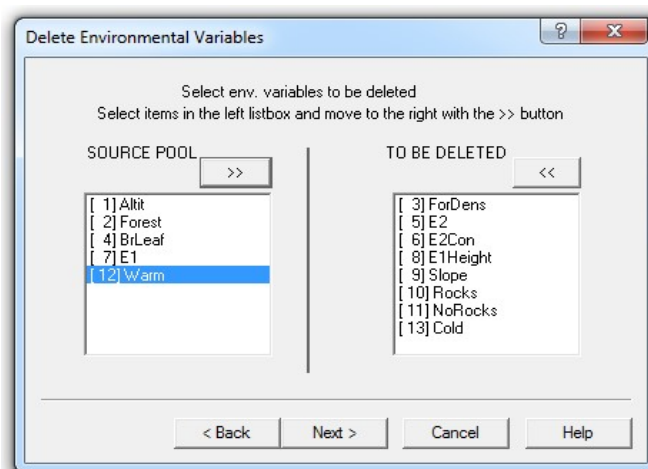
Marginal Effects			Conditional Effects				
Variable	Var.N	Lambda1	Variable	Var.N	LambdaA	P	F
Altit	1	0.11	Altit	1	0.11	0.002	5.32
BrLeaf	4	0.1	BrLeaf	4	0.08	0.002	3.72
Forest	2	0.09	Forest	2	0.06	0.002	2.91
E2Con	6	0.07	Warm	12	0.03	0.032	1.65
Slope	9	0.05	E1	7	0.03	0.038	1.71
E1	7	0.05	E2	5	0.03	0.066	1.54
Warm	12	0.05	E2Con	6	0.03	0.064	1.6
Cold	13	0.05	Slope	9	0.02	0.144	1.34
E1Height	8	0.05	ForDens	3	0.02	0.39	1.05
E2	5	0.04	E1Height	8	0.02	0.548	0.94
Rocks	10	0.04	Rocks	10	0.01	0.606	0.88
NoRocks	11	0.04					
ForDens	3	0.03					

Din acest tabel de rezultate vom vedea în coloana din stânga modul în care sunt ordonate variabilele de mediu în sensul descrescător al importanței lor (Lambda1, de sus în jos) în ceea ce privește explicarea variației în compoziția de specii, iar în dreapta vom vedea modul de grupare al variabilelor în funcție de îndepărtarea acumulativă a variației explicate de variabilele de mediu, în ordinea inversă a importanței lor. Când probabilitatea (P) a grupului de variabile în coloana de *Conditional Effects* depășește nivelul de asigurare ales, de obicei de 0.05, oprim modelul, respectiv selecția aditivă de variabile, la grupul astfel delimitat. Aceasta înseamnă că variabilele care explică cel mai bine variația în compoziția speciilor de păsări sunt, în ordine aditivă: Altit (altitudinea), la care se adaugă BrLeaf (frecvența relativă a foioaselor în arboret), Forest (acoperirea cu arbori), Warm (suprafețe expuse spre soare, sud și direcții asemănătoare) și E1 (acoperirea stratului ierbos). Prin adăugarea de alte variabile modelul nu mai are de câștigat, deoarece acestea sunt redundante. Prin urmare putem selecta numai aceste variabile, care alcătuiesc cel mai bun model de explicare a structurii comunităților de păsări și să reluăm graficul. Vom elimina celelalte variabile, prin întoarcerea la opțiunile anterioare: în fereastra de *Data Editing Choices* vom selecta *Delete Environmental Variables*, ca mai jos:





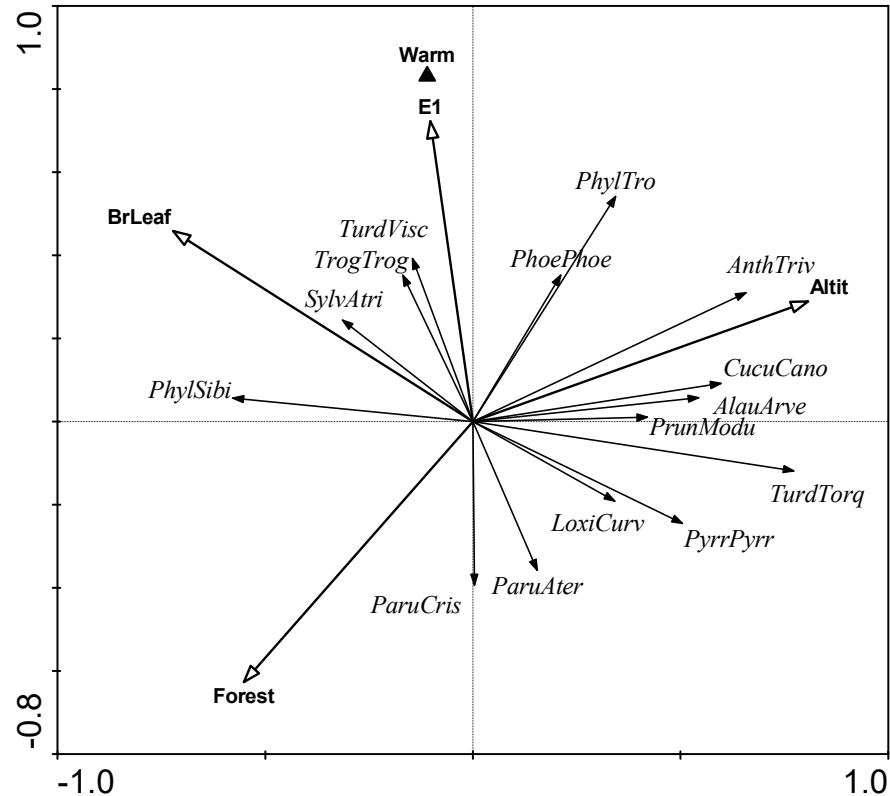
iar din următoarea fereastră vom selecta variabilele pe care dorim să le îndepărtăm, astfel:



\*\*\*\* Summary \*\*\*\*

Axes	1	2	3	4	Total variance
Eigenvalues	0.154	0.071	0.039	0.032	1.000
Species-environment correlations	0.877	0.774	0.709	0.715	
Cumulative percentage variance					
of species data	15.4	22.5	26.4	29.6	
of species-environment relation:	49.7	72.7	85.3	95.5	
Sum of all eigenvalues					1.000
Sum of all canonical eigenvalues					0.310

Rezultatele ne spun că dacă selectăm numai cele 5 variabile de mediu, acestea explică 31% din compoziția specifică, față de toate variabilele considerate, care adaugă numai cca. 13% în plus. În studiul viitor, sau la continuarea unor studii similare, ne-am putea decide astfel să evaluăm numai grupul celor 5 variabile alese, care realizează cel mai bun și semnificativ model. Graficul, după ce vom excluda speciile mai puțin reprezentate cantitativ, va arăta astfel:



### Analiza directă a gradientului: efectul altitudinii

Odată realizate aceste analize, putem diversifica întrebările și interpretările datelor noastre. Am aflat că altitudinea este foarte importantă în explicarea compoziției specifice a avifaunei. Am putea dori să aflăm cât de importantă este această variabilă, dacă ar fi singura care ar acționa și cum răspund comunitățile de păsări la celelalte variabile, dacă îndepărtăm varianța explicată de altitudine?

**Prima întrebare:** cum se explică compoziția comunităților de păsări dacă considerăm numai altitudinea medie a suprafețelor de eșantionare, ca singura variabilă independentă?

#### Mod de lucru:

- salvăm proiectul original sub un nume nou (de exemplu rda1.con), utilizând din meniul principal File → Save as...→rda1.con
- acum trebuie realizate niște modificări din *Options*. Ne asigurăm că în prima pagină (prima fereastră) alegem *extract patterns from explained variation only*, apoi alegem un nou nume pentru fișierul de soluții (de exemplu rda1.sol), la *Type of Analysis* vom selecta *RDA*, după care vom lăsa selecțiile de pe celelalte pagini așa cum sunt până la *Data Editing Choices*. Aici vom selecta opțiunea de ștergere a variabilelor de mediu

(*Delete Env. variables*), iar în următoarea fereastră (*Delete Environmental Variables*) vom selecta toate variabilele, cu excepția primei = altitudine (*Altit*) și le vom muta cu comanda (butonul) > > în coloana din dreapta ferestrei (*to be deleted*). În rest lăsăm selecțiile implicite și la fereastra de *Global Permutation Test* vom selecta oricare dintre teste, deoarece având o singură variabilă explicativă (altitudinea), vom avea implicit și un singur test pentru axa canonică unică astfel definită. În următoarele ferestre vom selecta *Unrestricted permutation* și *Finish*.

După ce vom realiza analiza (*Analyze*) rezultatul testului de permutări Monte Carlo va indica faptul că ipoteza nulă, care afirmă lipsa de relație între altitudine și compoziția specifică, este respinsă la  $p = 0.002$ , prin urmare avem o relație semnificativă.

\*\*\*\* Summary of Monte Carlo test \*\*\*\*

Test of significance of all canonical axes : Trace = 0.115  
 F-ratio = 5.323  
 P-value = 0.0020

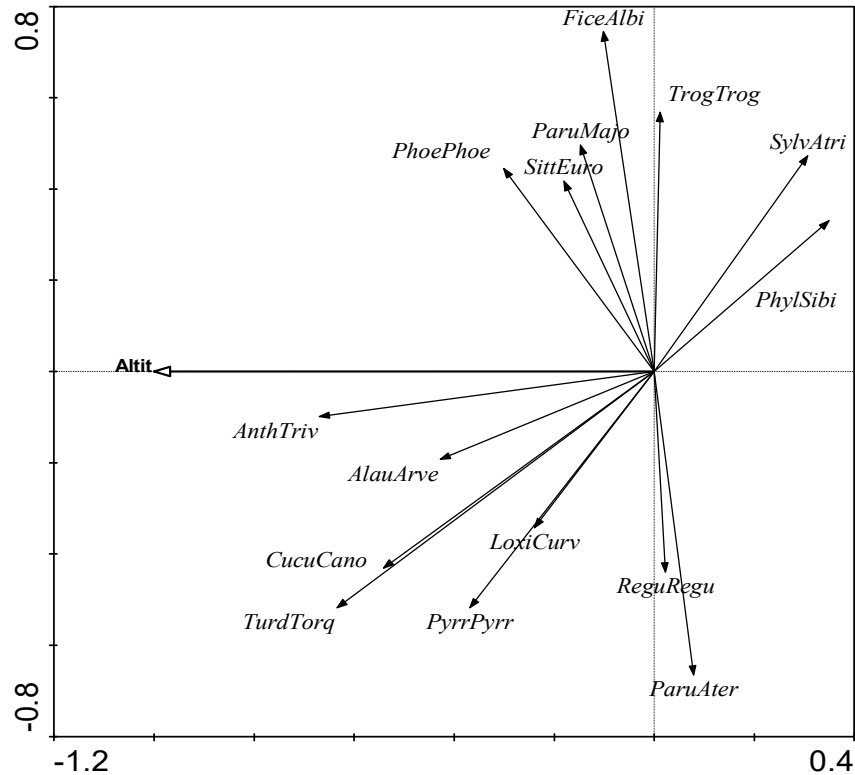
(499 permutations under reduced model)

Rezultatele analizei arată ca mai jos:

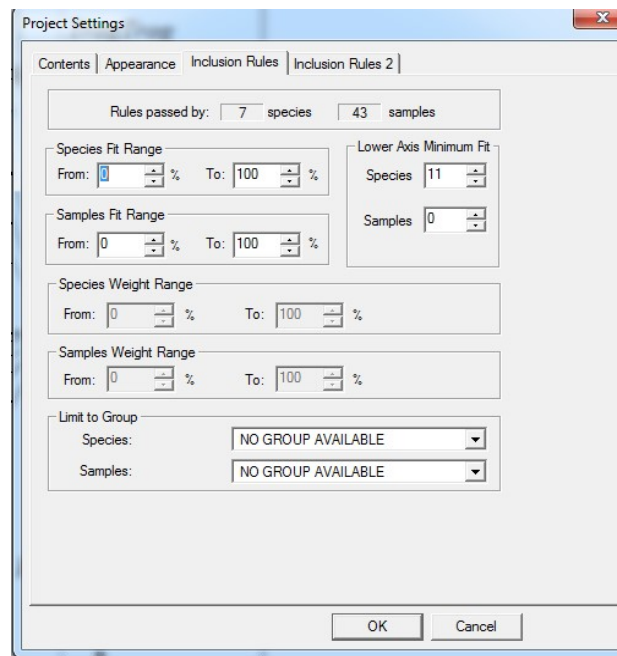
\*\*\*\* Summary \*\*\*\*

Axes	1	2	3	4	Total variance
Eigenvalues	0.115	0.152	0.118	0.075	1.000
Species-environment correlations	0.792	0.000	0.000	0.000	
Cumulative percentage variance					
of species data	11.5	26.7	38.5	46.0	
of species-environment relation:	100.0	0.0	0.0	0.0	
Sum of all eigenvalues					1.000
Sum of all canonical eigenvalues					0.115

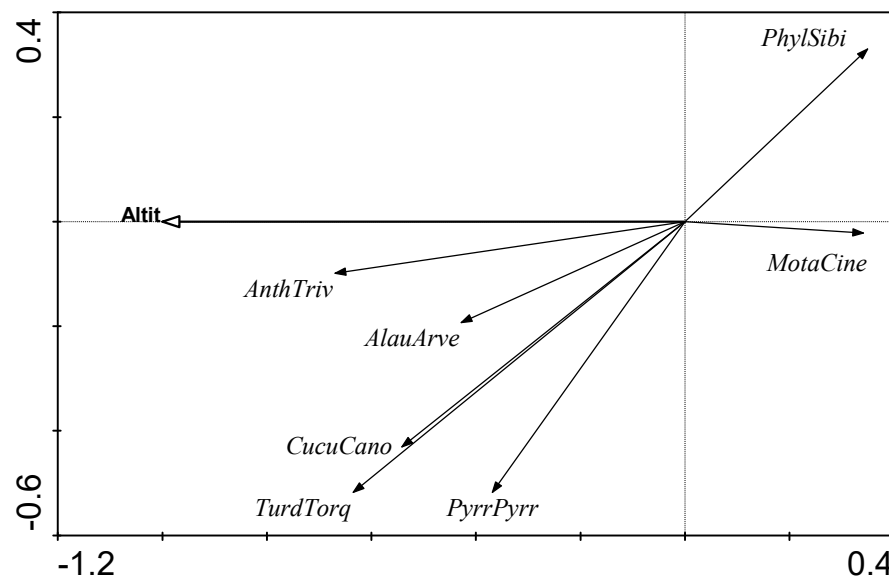
Înmulțind cu 100 valoarea rădăcinii latente a primei axe (deoarece există o singură variabilă, axa canonică și suma valorilor va fi aceeași) respectiv *sum of all canonical eigenvalues* =  $0.115 \cdot 100 = 11.5\%$ , înseamnă că altitudinea explică 11.5% din variația compoziției în specii a comunităților de păsări. Ne aducem aminte că toate variabilele explicau 44.2%, prin urmare altitudinea este răspunzătoare de peste un sfert din variația speciilor, dar următoarele două axe (2 și 3) care nu sunt canonice (neconstrânse) explică mai multă variație decât altitudinea (15.2% respectiv 11.8%).



În graficul de mai sus redăm biplotul specii-mediu din RDA care rezumă diferențele în compoziția comunităților de păsări de-a lungul gradientului altitudinal. Subliniem că acest grafic este de tip mixt (hibrid), în sensul că prima axă este canonică, în timp ce cea de-a doua nu. Se recomandă ca reprezentarea să includă numai speciile care sunt puternic corelate cu axa altitudinală (axa 1) iar celelalte să fie eliminate (adică toate cu unghiuri apropiate de  $90^0$  față de prima axă, sau apropiate de a doua). Deoarece altitudinea explică cca. 11% din variabilitatea datelor de specii, putem selecta numai speciile care sunt explicate în proporție de cel puțin 11% de către această axă. Această condiție sau limită se poate fixa în *CanoDraw* din fereastra de dialog *Project* → *Settings* → *Inclusion Rule* → după care vom selecta câmpul *Species* în *Lower Axis Minimum Fit*, deoarece numai axa inferioară (*lower*) este implicată. Numai șapte specii vor trece de această limită sau condiție.



Graficul este redat mai jos (în meniul *Create* se va alege opțiunea de *Recreate graph*)



### Analiză directă de gradient: efectul condițional al celorlalte caracteristici de habitat

Vom termina acest studiu de caz cu o analiză parțială canonică mai avansată, care va răspunde la întrebarea: **putem să detectăm vreun efect semnificativ al celorlalți descriptori de habitat, dacă îndepărtăm variabilitatea compozițională explicată de către altitudinea medie a pătratului de probă?**

La această întrebare vom răspunde printr-o analiză RDA din care altitudinea va fi considerată ca și **covariată**.

Mod de lucru:

- *File* → *Save as* (alegem un nou nume de proiect, de exemplu rda2.con) → *clear the project log*;

- *Options* - iar din prima fereastră (*Available Data*) vom selecta opțiunea a treia: *Species, environmental and covariable data available* → *Next*

- În fereastra de definire a datelor (*Data Files*) va trebui să identificăm fișierul din care programul va alege covariatele. Deoarece acestea sunt tot în fișierul birds\_env.dta, vom alege acest fișier, sau vom copia conținutul câmpului *Environment data file name* la *Covariables data file name*.

Urmează *Next* → *RDA* → etc. neschimbat până la *Data Editing Choices*. Programul deja a selectat opțiunea de ștergere (*Delete*) a covariatelor (*Covariables*), iar după *Next*, vom selecta variabilele corespunzătoare.

- În fereastra *Delete Environmental Variables* vom șterge (muta în dreapta) numai altitudinea (*Altit*), în timp ce toate celelalte vor fi în stânga (variabilele 2 - 13).

- *Next* → *Delete covariables*: în această fereastră vom face exact invers. Ștergem (mutăm în dreapta) toate celelalte variabile exceptând altitudinea (variabila 1). Numai variabila *Altit* va rămâne în stânga. În următoarea fereastră (după *Next*) vom lăsa opțiunea selectată *Do not use forward selection*.

- În fereastra *Global Permutation Test* va fi deja selectată opțiunea *Significance of canonical axes together* (lăsăm ca atare), iar în următoarea lăsăm opțiunea *Unrestricted permutations* selectată. În final închidem proiectul pe următoarea pagină (*Finish*). Soicimăm analiza (*Analyze*).

- În fereastra de rezultate sau în jurnal (*Log View*), vom obține, alături de alte rezultate, următoarele valori:

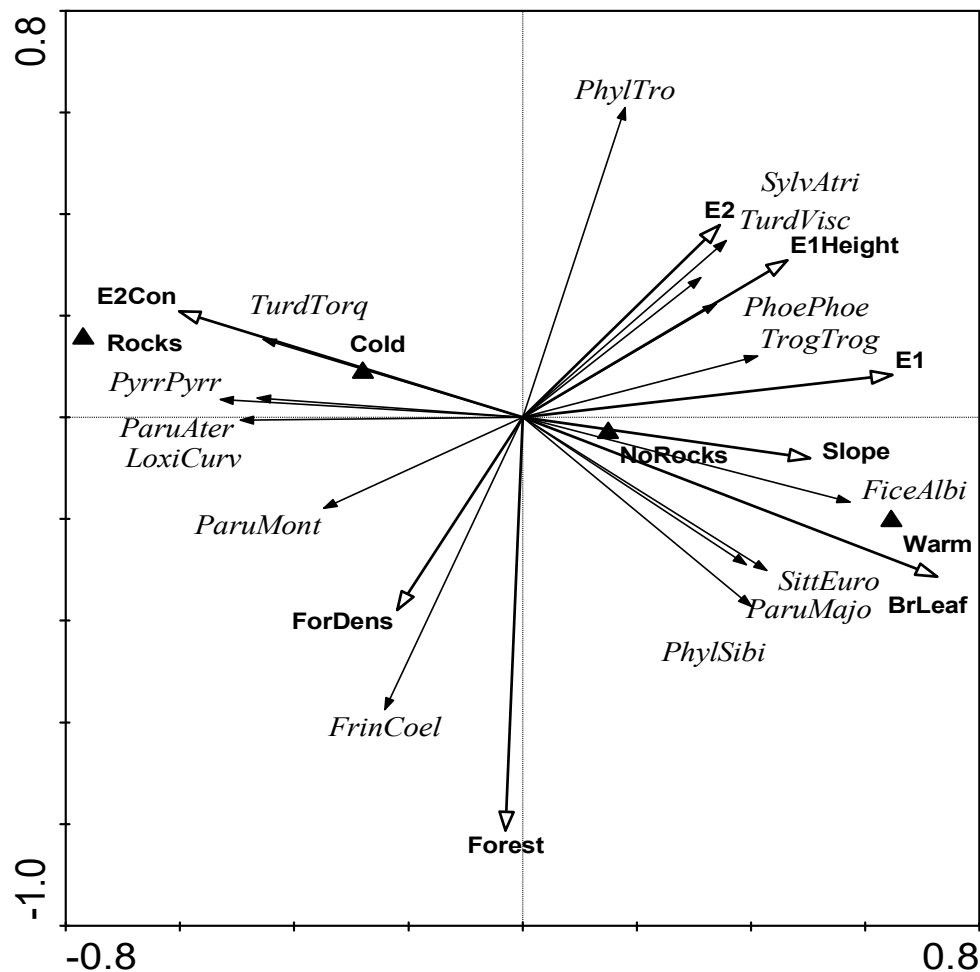
\*\*\*\* Summary \*\*\*\*

Axes	1	2	3	4	Total variance
Eigenvalues	0.112	0.067	0.042	0.034	1.000
Species-environment correlations	0.886	0.810	0.818	0.765	
Cumulative percentage variance					
of species data	12.6	20.3	25.0	28.8	
of species-environment relation:	34.2	54.7	67.4	77.8	
Sum of all eigenvalues					0.885
Sum of all canonical eigenvalues					0.328

The sum of all eigenvalues is after fitting covariables  
Percentages are taken with respect to residual variances  
i.e. variances after fitting covariables

All four eigenvalues reported above are canonical and correspond to axes that are constrained by the environmental variables.

Se poate observa că variabilitatea explicată de ceilalți descriptori de habitat, în adăuție față de informația adusă de altitudine, este destul de mare (32.8% din totalul variabilității în compoziția specifică a comunităților de păsări), prin comparație cu variația explicată de altitudine (11.5%, vizibilă din acest tabel și ca diferența *Total variance - Sum of all eigenvalues* = 1.000 - 0.885). Dar este evident că modelul de față este mult mai complex ca număr de grade de libertate. Acest tabel de rezultate este urmat de raportul testului de permutări Monte-Carlo. Se poate vedea că contribuția adăuțională a celor 12 descriptori este extrem de semnificativă ( $p = 0.002$ ).



*RDA de tip biplot care redă efectul descriptorilor de habitat (al variabilelor de mediu) asupra compoziției comunităților de păsări, după ce a fost îndepărtat efectul gradientului altitudinal (altitudinea este considerată covariată).*

Diagrama biplot specii-mediu care rezumă efectul adăuțional al celorlalte caracteristici de habitat, după ce efectul altitudinii a fost luat în considerare, este redat mai jos. Nu uităm să declarăm variabilele nominale (*Project* → *Nominal variables* →

*Environmental variables* → selecție variabile substitutive sau nominale). Numai speciile cele mai bine reprezentate sau cu ponderi mai ridicate (*Project* → *Settings* → *Species Fit Range > 15*) au fost selectate aici, pentru a facilita interpretarea graficului.

Observăm că prima axă a acestui RDA reflectă parțial diferențele dintre pătratele dominate de fag (pe partea dreaptă a graficului) și arboretele dominate de molid (pe partea stângă a diagramei), precum și dezvoltarea stratului ierbos (mai dezvoltat în jumătatea dreaptă a figurii). A doua axă de ordonare se corelează negativ cu acoperirea totală cu arbori (*Forest*) a pătratelor de probă și, într-o măsură mai mică, cu densitatea medie a pădurii (*ForDens*).