

10 AVR. 1972

1

Distances non-paramétriques
entre profils

Raymond Van den Driessche et Ana García Gómez

† Ir. Responsable de la Banque de Données Pédologiques de l'ORSTOM, 70, route d'Aulnay, F-93-Bondy.
Dr Sc. 6, rue Rabelais, F-93-Bondy.

14 AVR. 1972

O. R. S. T. O. M.

Collection de Référence

n° 5384 Biom

Summary

Soil scientists familiar with storage and retrieval of field data are aware of the need to manipulate simultaneously ordinal and interval variables.

A distribution-free measure of distances between profiles using ranks is exemplified. The distances are measured on the basis of an incomplete matrix of m rows or profiles and v columns or variables. The column entries are the same variables taken in turn in all the horizons named A1, A3, B, BC, C, C1, C2. In the numerical example from São Tomé isl. Afr. the matrix 28 x 98 consists of 28 profiles and 14 variables (hue, value, chroma, clay, silt, fine sand, coarse sand, organic matter, C/N, pH, exchangeable calcium, sum of the cations, cation exchange capacity, saturation level) in several of the 7 horizons. Therefore, the data matrix is only half filled with data, and is not restricted to field data.

Punched output in SF9.4 of the lower half of the distance matrix is followed by a printed symbolic output in condensed form so as to allow for $m \leq 127$ ($v \leq 2000$ and $mv \leq 93000$) on the 1108 computer.

Introduction

Lors du traitement de données descriptives des horizons ou de l'environnement du profil, nous ressentons le besoin de disposer d'une mesure de distance non-paramétrique. Les hypothèses d'homoscédasticité et de normalité faites, si souvent, pour les variables de laboratoire ne sont plus de mise dans ce domaine par excellence des variables qualitatives. Bien que la méthode de rangs que nous présentons pour cette mesure n'ait pas encore fait, à notre connaissance, l'objet d'applications, nous la couplons à une mesure originale de distance entre profils entiers. Les définitions indispensables sont rappelées et la méthode est appliquée à un lot de 28 profils de St Tomé.

Rappel de définitions

Le GROUPE est un ensemble de données observées de différente nature. Il y a une seule donnée par variable. Constituent le plus souvent un groupe, toutes les données observées dans un des horizons d'un profil. Peut aussi former un groupe, la totalité des données observées dans un ensemble de profils classés de la même façon dans la classification des sols. Nous sommes d'avis de former le groupe avec toutes les données observées dans la totalité des horizons d'un même profil, mais en respectant l'identité de chaque horizon.

Opposée au groupe, la VARIABLE est un ensemble de données observables, de même nature, exprimées dans la même unité ou vocable normalisé, mais dont une seule est observée dans un même groupe. La variable "test de plasticité" est utilisée quand l'horizon est humide ou très humide; elle se matérialise par l'aptitude que possède le matériau à subir un changement continu de forme. Pour cette variable, prise à titre d'exemple, quatre données sont observables : non plastique, peu plastique, plastique, très plastique. Second exemple tiré de la description des horizons, la variable "nature des croûtes et/ou efflorescences éventuelles" offre à l'observateur quatre données : chlorurées, bicarbonatées, carbonatées, sulfatées, dont une seule est observée.

Une DONNÉE, et seulement une, exprime donc l'état d'une variable dans un groupe déterminé. Les données sont exprimées en langage naturel ou en code.

Le MANQUANT est une donnée qui fait défaut et cela pour des motifs liés aussi bien à la nature même des données qu'à l'observateur.

Les IDENTIQUES sont les données d'une même variable qui ont la même expression, en langage naturel ou en code, dans deux ou plusieurs groupes. Des données ayant la même expression mais dépendant de variables différentes ne sont pas des identiques.

Le TABLEAU est un assemblage ordonné de groupes, de variables, de données et de manquants. Il y a autant de lignes que de groupes, autant de colonnes que de variables.

Exposé de la méthode

Considérons le TABLEAU ci-après. Il est constitué de données en provenance de St Tomé et nous servira d'application numérique.

L'île de St Tomé est située entre les longitudes E 6-45-0 et 6-28-19 et entre les latitudes N 0-24-30 et 0-0-0. Elle a été étudiée par les pédologues Cardoso & Garcia (1962).

Les auteurs ont donné l'identité des horizons décrits, en utilisant entre autres les symboles A1, A3, B, BC, C, C1, C2. Cela nous a permis d'extraire de leur mémoire les 28 profils n° 44A, 93, 147, 204, 125, 33A, 82, 88, 157, 212, 209, 41, 38, 111, 131, 6, 47, 56, 28A, 29, 45, 207, 23, 36A, 37, 94, 100, 69. Ces profils sont représentatifs des sols "paraferralíticos" et "fersialíticos tropicaux" de l'île. L'application numérique porte donc sur 28 groupes.

Les variables retenues par nous sont, dans l'ordre, la teinte, la valeur et l'intensité Munsell de l'horizon, l'argile, le limon, le sable fin, le sable grossier, la matière organique, le rapport C/N, le pH, le calcium échangeable, la somme des cations échangeables, la capacité d'échange, le taux de saturation. L'originalité de notre approche réside dans le fait de considérer une variable déterminée (limon par exemple) dans sept horizons différents (A1, A3, B, BC, C, C1, C2) comme autant de variables différentes. Tout reste individualisé dans notre tableau : A1 et A3 ; B et BC ou BC et C ; C, C1 et C2 ne sont pas confondus.

Le GROUPE correspondant à un profil est, par conséquent, un ensemble de DONNEES provenant des VARIABLES des horizons décrits et de MANQUANTS. Dans cette optique, les manquants correspondent, à la fois, aux variables non décrites dans un horizon quelconque et aux variables des horizons n'existant pas dans le profil qui est décrit mais bien dans d'autres profils de l'île. Nous désignons les manquants par le code -1. L'effectif d'une variable est le nombre de DONNEES de cette variable sur l'ensemble des groupes. Si m désigne le nombre de groupes, m_i est l'effectif de la variable i. L'effectif de la première variable est 28 ; celui de la dernière est 3.

Il nous reste à souligner quelques identiques dans ce tableau. De la variable 1, la donnée 15 est présente cinq fois, la donnée 17,5 deux fois, la donnée 20 vingt fois; ce qui fait trois lots d'identiques.

Les RANGS r_{ki} sont des numéros d'ordre 1 à m_i affectés aux m_i données x_{ki} d'une variable i . Pour les e_i lots d'identiques de la variable i l'affectation des rangs suit la règle suivante : les t_{qi} identiques du même lot q sont remplacés par la moyenne arithmétique des t_{qi} rangs suivants non encore utilisés.

L'égalité entre le rang maximal et l'effectif n'est vérifiée que lorsque la donnée maximale est seule. Pour une donnée maximale en t_{qi} identiques, le rang maximal est égal à $m_i - (t_{qi} / 2) + 0,5$

Premier exemple :

Les données de la variable 98 = taux de saturation dans l'horizon C2 sont remplacées par les rangs

k	$x_{k 98}$	$r_{k 98}$
3	12,3	1
4	38,8	3
15	15,9	2

$e_{98} = 0$ il n'y a pas d'identiques

$m_{98} = 3$ l'effectif est de 3 données

Deuxième exemple :

Pour la variable 43 = teinte Munsell de l'horizon BC nous remplaçons

dans les groupes k les données $x_{k 43}$ par les rangs $r_{k 43}$

1	15	2,5
2	15	2,5
6	20	10
7	20	10
8	20	10
9	20	10
11	17,5	4
12	20	10
13	20	10
16	20	10
18	12,5	1
20	20	10
21	20	10
25	20	10
27	20	10

L'effectif $m_{43} = 15$. La donnée minimale 12,5 est remplacée par le rang 1. Les deux identiques 15 du lot 1 sont remplacés par la moyenne arithmétique des deux rangs non encore utilisés $(2+3)/2 = 2,5$. La donnée suivante 17,5 est présente une seule fois, le rang 4 la remplace. Les onze identiques 20 sont remplacés par la moyenne arithmétique des onze rangs non encore utilisés, c'est-à-dire par

$$(5+6+7+8+9+10+11+12+13+14+15)/11=10$$

Pour mémoire, $e_{43} = 2$; $t_{1 43} = 2$; $t_{2 43} = 11$.

La DISTANCE est une mesure de dissemblance entre groupes. Elle est comprise entre 0 et 1, est indépendante de la fonction de répartition des variables et est proche de celle de Kendall & Stuart (1966).

$$D_{hk} = \frac{1}{v} \sum_{i=1}^v \frac{(r_{hi} - r_{ki})^2}{I_i}$$

mesure la distance entre le groupe h et le groupe k .

La correction

$$I_i = \left[m_i^3 - m_i - \sum_{q=1}^{e_i} (t_{qi}^3 - t_{qi}) \right] / 12$$

est calculée pour chaque variable à partir de :

- m_i effectif de la variable i
- e_i nombre de lots d'identiques pour la variable i
- t_{qi} nombre d'identiques dans le lot q de la variable i.

Pour toute distance calculée entre deux groupes, c'est pour chacune des variables le carré de la différence entre rangs homologues, divisé par la correction, qui seul intervient. La correction prend en compte, à la fois, l'effectif de la variable, le nombre de lots d'identiques et le nombre d'identiques par lot. Elle est nulle et la variable abandonnée en conséquence lorsque $t_{qi} = m_i$ c'est-à-dire quand les données sont toutes identiques ou quand il n'y a qu'une seule donnée.

Résultats de l'application numérique

Le programme KAV écrit pour nous par A.M. Aubry de la Banque de Données Pédologiques et transposé sur 1108 par F. Savary de STAD a permis de calculer en 12 secondes les 378 distances que voici. C'est dans ce format SF9.4 que la moitié inférieure de la matrice d'ordre 28 est disponible pour un traitement ultérieur.

Profils	1	2	3	4	5	6	7	8
	9	10	11	12	13	14	15	16
	17	18	19	20	21	22	23	24
	25	26	27					
2	.0844							
3	.0478	.0772						
4	.0371	.0435	.4425					
5	.1115	.0857	.1324	.0926				
6	.0950	.1110	.0449	.0595	.1177			
7	.1053	.1370	.0759	.1147	.1387	.1096		
8	.0964	.1013	.0666	.0791	.1309	.0574	.0761	
9	.0608	.0814	.0530	.0763	.1608	.0907	.0716	.0997
10	.0733	.0598	.0895	.0691	.1170	.0817	.1014	.0573
10	.0943							
11	.0583	.0911	.1096	.0547	.0877	.0960	.0966	.0921
11	.1053	.0663						
12	.1355	.1128	.1348	.1296	.0658	.0903	.1222	.0834
12	.1397	.1247	.1146					
13	.0895	.0542	.0807	.0919	.1103	.0813	.1126	.1046
13	.0563	.0781	.0940	.0946				
14	.0579	.0493	.0602	.0661	.1524	.0423	.1151	.0714
14	.0407	.1064	.0872	.1076	.0165			
15	.2293	.0866	.3263	.6909	.1363	.0703	.0661	.0621
15	.0539	.0874	.1066	.0455	.0451	.0506		
16	.1091	.0727	.0892	.0869	.1125	.0999	.1382	.1313
16	.0883	.1092	.0883	.1286	.0553	.0627	.0702	
17	.1115	.0790	.0944	.0441	.1624	.0947	.2235	.1603
17	.1479	.1417	.1002	.1915	.1294	.1211	.1465	.0964
18	.1048	.0669	.0582	.0449	.0790	.1119	.1618	.1338
18	.1139	.1271	.1192	.1483	.1338	.0876	.1032	.0910
18	.0953							
19	.0922	.0720	.0942	.0670	.1510	.0801	.1554	.1068
19	.1014	.1139	.0897	.1448	.0895	.0737	.1010	.0707
19	.0382	.0962						
20	.0915	.0542	.0705	.0688	.1136	.0812	.0865	.0985
20	.0727	.0996	.0735	.1181	.0540	.0795	.0344	.0330
20	.1448	.0730	.0814					
21	.0863	.0612	.0793	.0694	.1536	.0670	.1232	.0803
21	.0510	.0819	.1092	.1090	.0585	.0490	.0829	.0711
21	.0944	.0990	.0502	.0453				
22	.0796	.0578	.1110	.0615	.1368	.0709	.1634	.0899
22	.0950	.0573	.0543	.1166	.0819	.0909	.1161	.0773
22	.0863	.1007	.0690	.0860	.0438			
23	.0964	.0652	.1146	.0730	.1974	.1092	.1951	.1372
23	.0912	.1285	.1022	.1817	.1017	.0832	.1077	.0839
23	.0600	.1129	.0249	.0906	.0418	.0766		
24	.0945	.0843	.1071	.0675	.0668	.0627	.1523	.1183
24	.1516	.0954	.0793	.0950	.1004	.0973	.1501	.1024
24	.0466	.0916	.0466	.0765	.1191	.0700	.0853	
25	.1451	.0929	.1197	.0792	.0774	.1103	.1907	.1213
25	.1498	.0659	.1077	.0829	.1067	.0857	.1134	.1139
25	.0405	.1401	.0340	.1081	.0990	.0299	.0509	.0384
26	.0731	.0854	.1079	.0564	.1375	.0989	.1465	.1042
26	.0978	.0974	.0799	.1377	.0995	.1006	.1331	.0871
26	.0921	.1069	.0501	.0713	.0471	.0557	.0399	.1156
26	.0412							
27	.1441	.1199	.1253	.0977	.1591	.1322	.2169	.1168
27	.1476	.1167	.1266	.1405	.1257	.0991	.0817	.1153
27	.0796	.1787	.0495	.1530	.1140	.1037	.0521	.0977
27	.0450	.0447						
28	.1004	.0603	.0677	.0746	.1328	.0443	.0829	.0670
28	.0758	.1070	.0608	.1200	.0372	.0373	.0429	.0290
28	.1227	.0850	.0776	.0536	.0746	.0948	.0993	.0621
28	.0479	.1095	.1237					

Pour en faciliter la lecture, les mêmes distances sont remplacées, ci-après, par des symboles de classe : 0, 1, 2, ... 9, A, B, ... N, P, ... Z. Les distances dépassant 0,3400 sont ainsi abandonnées.

```

      .
      .      111111111222222222
      .1234567890123456789012345678
      .....
1..954CABA786E96NBCDAA98AAF8FB
2..9.859CEB96AC659887867679A9C7
3..58. E58769BE97Y9A6A88CCBCBD7
4..45 .A6C8876DA7 95577778786A8
5..C9EA.CEEHC97CGECH8GCGEK78EGE
6..AC56C.B6A9AA958AAC997887CAE5
7..BE8CEB.88BADCC7ENH9D9D9GKFM9
8..AB78E68.A6A9B87EHEBA99ECDBC7
9..7968HA8A.ABE6569FCB86AA6FAF8
10..8697C9B6A.7D8B9BFDCA96DA7ACB
11..6AB69AAAB7.CA9B9BC9BB6B8B8D7
12..ECED7AD9EDC.AB5DKFFCBCJA9EFC
13..969AC9CB68AA.256DE9669BBBBD4
14..657765C85B9B2.67D9885A9A9BA4
15..N9Y E87769B556.8FBB49CBGCE95
16..B899CAEE9B9D678.AA84889BC9C3
17..C8A5HANHFFBKDDFA.A4FA9755ABD
18..B7658CHECDCFE9BAA.A8ABCAFB19
19..ABA769GBBC9F98B84A.967354658
20..A687C99A8A8C6844F89.59A8B8G6
21..978767D969BB6598AA65.55CA5C8
22..86C7EBH9A66C9AC89B795.8836BA
23..A7C8KBKEADBUB9B97C3A58.9646A
24..A9B777GCGA8ABAGB5A58C89.4CA7
25..FAC88CKDF7B9B9CC5F4BA364.555
26..89B6EAFBAA8EABE9AB68564C5.5B
27..FCDAGEMCFCDFA9C8I5GC86A55.D
28..B778E5978B7C4453D9868AA75BD.

```

Présentée sous cette forme matricielle, la distance est lue à l'intersection de la ligne et de la colonne qui portent, respectivement, les numéros des deux profils entre lesquels la distance est mesurée.

A titre d'exemple, les résultats sont lus ainsi : entre les profils 13 et 14 (ligne 13, colonne 14) la distance est 2 (classe 0,0101 à 0,0200), entre les profils 3 et 6 (intersection de la ligne 3 et de la colonne 6) la distance est 5 (classe 0,0401 à 0,0500), entre les profils 3 et 11 le code B représente une distance de la classe 0,1001 à 0,1100. La distance étant une mesure relative, comprise entre 0 et 1, et le tiers inférieur méritant une attention plus particulière, la représentation ci-dessus semble suffire.

Bibliographie

Cardoso J. Carvalho, Garcia J. Sacadura, 1962 - Carta dos solos de
São Tomé e Príncipe. Mem. Junta Invest. Ultram., 2a sér., nº 39,
Lisboa, 306 p.

Kendall M.G., Stuart A., 1966 - The advanced theory of statistics. vol. 3.
Design and analysis, and time series. Griffin, London, 552 p.