

Commentaires sur l'article de G. Celeux et C. Robert
G rard Govaert
CNRS URA 817 - Universit  de Technologie de Compi gne

Le texte de Celeux et Robert sur les effets parfois n fastes du recodage des variables continues en variables qualitatives appellent deux remarques de nature diff rente : la premi re concerne l'analyse des variables initiales qui aurait pu  tre un peu simplifi e et la seconde porte sur une analogie que l'on peut faire entre la discr tisation des images et le codage de variables continues en variables qualitatives.

0.1 L'analyse des variables continues

L'analyse pr alable des variables continues, telle qu'elle est pr conis e dans la conclusion du texte de Celeux et Robert peut  tre effectu e plus simplement et plus rapidement, sans faire appel notamment   l'analyse en composantes principales et en utilisant uniquement l'analyse des statistiques descriptives  l mentaires (histogrammes, tableau de corr lations, ...):

- l'histogramme de la variable z_4 met en  vidence la pr sence du papillon aberrant,
- la corr lation tr s forte de 0.97 entre les variables z_2 et z_3 permet d' liminer sans une grande perte d'information l'une des deux variables,
- enfin, la prise en compte de l'effet taille, assez classique pour ce type de donn es, conduit   ne retenir que les deux nouvelles variables $t_1 = z_1/z_4$ et $t_2 = z_2/z_4$.

Il suffit alors de repr senter le nuage des papillons dans le plan (t_1, t_2) (figure 1) pour mettre en  vidence tr s clairement la structure en trois classes de l'ensemble des papillons.

Remarquons que m me une simple analyse des diff rents plans (z_i, z_j) des variables initiales aurait pu suffire : en effet, le plan (z_3, z_4) (figure 2) montre assez nettement une structure en trois classes ainsi que le papillon aberrant.

0.2 Discr tisation des images et recodage des variables continues

En ce qui concerne la perte d'information due au recodage des variables continues en variables qualitatives, il est int ressant de faire un parall le avec le traitement d'images : si nous nous limitons, par exemple aux deux variables z_1 et z_2 , le nuage des papillons peut  tre consid r  comme une image dans le plan (z_1, z_2) . Comme la variable z_1 cod e en entier de 21   31 peut prendre 11 valeurs diff rentes et z_2 cod e en entier de 30   42 peut prendre 13 valeurs diff rentes, l'image ainsi d finie est constitu e de $11 \times 13 = 143$ pixels ou unit s  l mentaires d'image.

Avec cette analogie, le recodage des variables z_1, z_2 en variables qualitatives   3 modalit s s'interpr te alors comme une discr tisation de l'image qui passe de 143 pixels   9 pixels.

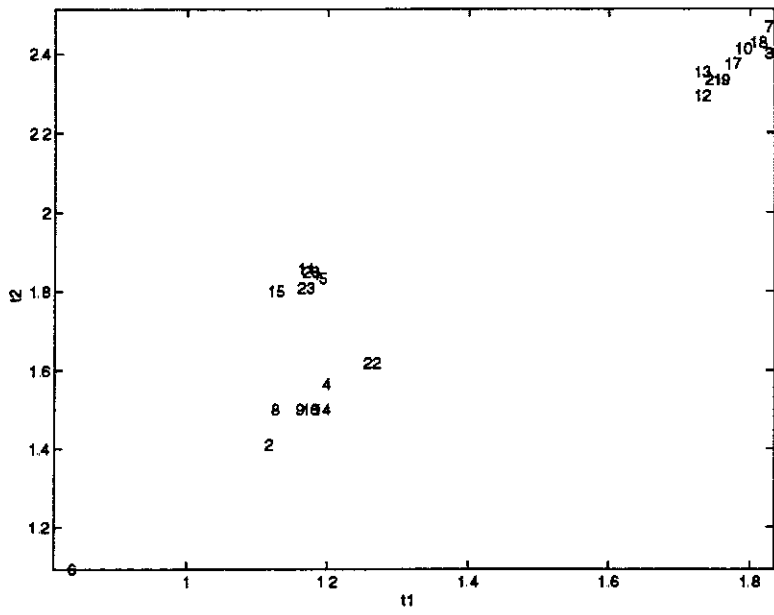


FIG. 1 - Plan t_1, t_2

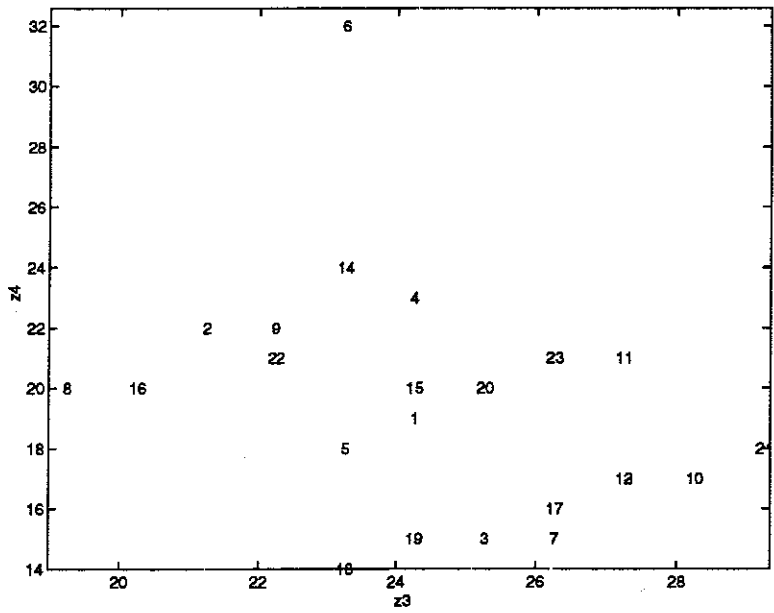


FIG. 2 - Plan z_3, z_4

Cette comparaison montre alors assez bien que ce qui paraît assez naturel et classique en analyse des données l'est beaucoup moins en traitement d'images. Il paraît en effet peu probable que s'il existait une forme intéressante et visible dans l'image initiale, elle soit encore visible dans l'image discrétisée.

Remarquons que pour notre problème complet, la diminution des pixels est encore plus importante: en effet, nous passons de $11 \times 13 \times 11 \times 18 = 28314$ unités élémentaires à $3 \times 3 \times 3 \times 3 = 81$ unités élémentaires. Dans ces conditions, pour que l'information pertinente soit conservée après une telle compression de l'information, il aurait fallu, comme le conclue le texte de Celeux et Robert, pour le moins que la discrétisation se soit faite en fonction de cette information (par exemple recherche d'une bonne rotation de l'image, ...), c'est-à-dire pour nous en prenant par exemple les meilleures combinaisons linéaires des variables avant d'en faire le recodage.