

LES ANALYSES EN CORRESPONDANCES MULTIPLES DANS SAS® ENTERPRISE GUIDE® 4.3

Le public francophone demande souvent s'il possible de faire une **ACM** (**A**nalyse en **C**orrespondances **M**ultiples) avec SAS Enterprise Guide. Ce type d'analyse n'existe pas en mode standard, dans une tâche. Néanmoins, une tâche expérimentale et téléchargeable sur notre site Internet peut être installée dans Enterprise Guide 4.3. Elle permet de réaliser ce type d'étude statistique très facilement en quelques clics.

Caractéristiques :

Catégorie : Enterprise Guide
OS : Windows
Version : 4.3
Vérifié en mars 2012

Nous allons débiter par une présentation rapide du concept de l'ACM. Ensuite nous expliquerons comment installer l'add-in pour réaliser ce type d'analyse dans Enterprise Guide. Enfin, nous illustrerons son fonctionnement au travers de deux exemples d'applications, et nous étudierons le code SAS qui s'exécute en arrière-plan.

Qu'est-ce qu'une ACM ?.....	1
Installation de l'Add-in	2
Exemples d'application de l'Add-in d'Analyse en Correspondances Multiples	3
Table avec des variables de classes.....	4
Application avec une table de contingence	6
Le code SAS traditionnel.....	8
Conclusion	9

Qu'est-ce qu'une ACM ?

L'ACM est une technique d'analyse de données, faisant partie de la famille de l'AFC (Analyse Factorielle des Correspondances), qui s'appuie sur des données catégorielles nominales. Les données peuvent se présenter sous forme d'un tableau de contingence : des variables en lignes, des variables en colonnes et les effectifs correspondants dans le tableau (par exemple, CSP versus lieu de vacances, ou encore lieu de résidence versus décennie) ou d'un tableau disjonctif complet (tableau de Burt), entre autres. Le premier type de tableau sera retenu pour l'étude qui va suivre. Généralement, les modalités dont le poids est inférieur à 5% sont regroupées avec d'autres (significativité statistique). Les objectifs de l'ACM sont de mettre en évidence:

- les relations entre les différentes modalités des variables,
- hypothétiquement, les relations entre les individus statistiques.

Pour illustrer la partie plus théorique qui va suivre, cette section s'appuie sur le cours de Thierry FOUCART (maître de conférences à la retraite de l'Université de Poitiers) qui étudie le comportement des fumeurs par rapport à leur sexe et leur majorité.

Exemple de tableau de contingence:

X : 3 modalités		Y : 4 modalités			
brunes : fumeur de brunes		mm : mineur masculin			
blondes : fumeur de blondes		mf : mineur féminin			
non fumeur		MF majeur féminin			
		MM : majeur masculin			
X	Y				
	mm	mf	MF	MM	
	brunes	63	37	41	47
	blondes	36	55	39	38
	non fumeur	34	27	72	38

Dans ce type d'analyse, on cherche les différences entre profil de lignes (ou colonnes, valable pour toutes les références suivantes au profil de lignes) plutôt que des différences individuelles.

Un profil de ligne est la répartition suivant les modalités en colonnes des individus données en lignes.

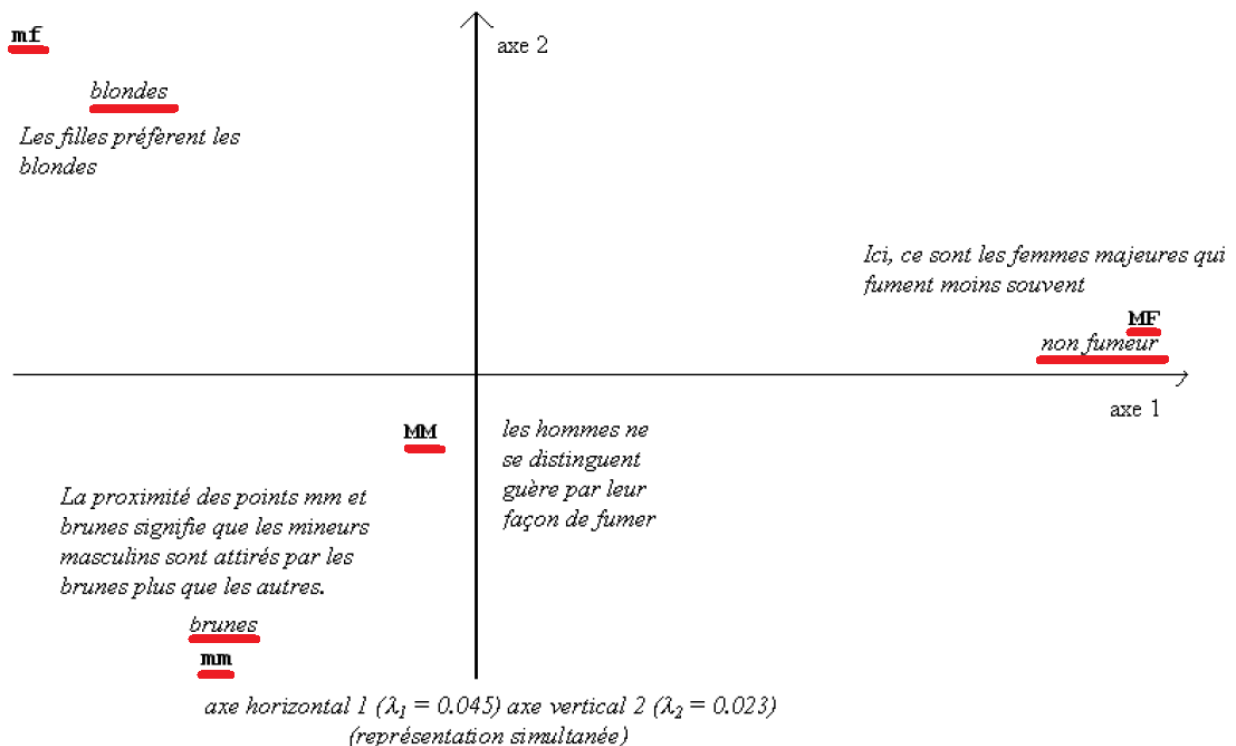
	mm	mf	MF	MM	total
profil brunes	0.335	0.197	0.218	0.250	1
profil blondes	0.214	0.327	0.232	0.226	1
profil non fumeur	0.199	0.158	0.421	0.222	1
centre de gravité P_j	0.252	0.226	0.288	0.233	1

profils lignes P_j^i

Une notion importante est P_j , par exemple 0,252 pour « mm » qui représente la part des mineurs mâles dans l'échantillon, ce qui correspond aux profils colonnes (sans tenir compte de leur comportement vis-à-vis de la cigarette). Comme il a été dit plus haut, le calcul des différences se fait entre les profils. Ainsi, la différence entre deux profils de lignes s'écrit :

$$d^2(i,i') = \sum_{j=1}^q [p_j^i - p_j^{i'}]^2 / p_{\cdot j}$$

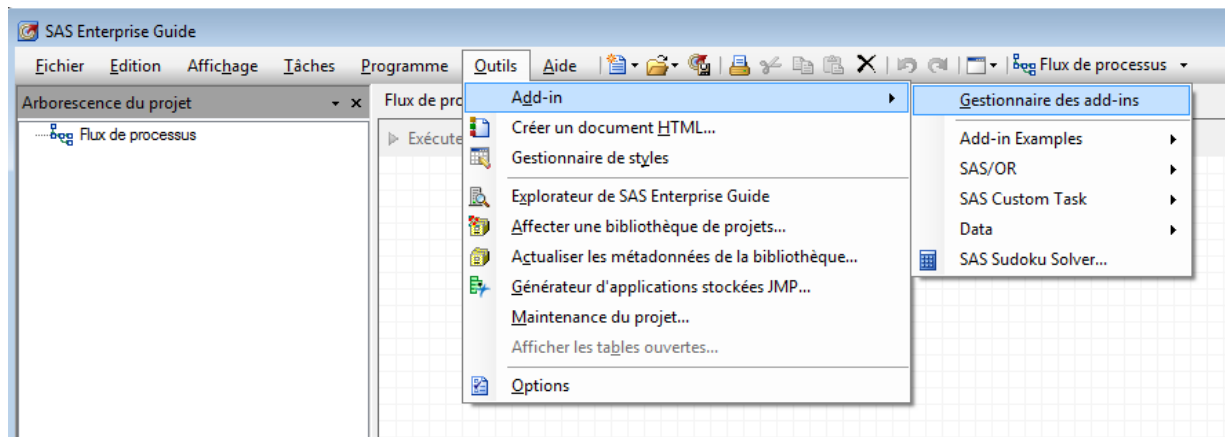
Les associations entre variables vont ensuite être découvertes en calculant la distance du chi-deux (présentée ci-dessus). Ces associations sont représentées graphiquement pour faciliter l'interprétation des structures dans les données. Les oppositions entre lignes et colonnes sont ensuite maximisées, pour découvrir les dimensions sous-jacentes les plus aptes à décrire les oppositions centrales dans les données. Comme en AFC ou en ACP, les axes d'analyses sont classés par ordre d'importance (en fonction de la quantité de variance, inertie du nuage, expliquée). Pour déterminer le nombre d'axes qu'il est pertinent d'étudier, plusieurs règles existent dans la littérature (100%/nombre d'axes=seuil de pertinence d'un axe). Une fois les axes pertinents choisis, il ne reste plus qu'à les analyser, au travers de différents critères (contributions, cosinus,...) qui ne seront pas détaillés dans ce papier. Voici le type de graphique, et quelques conclusions, qui sont obtenus :



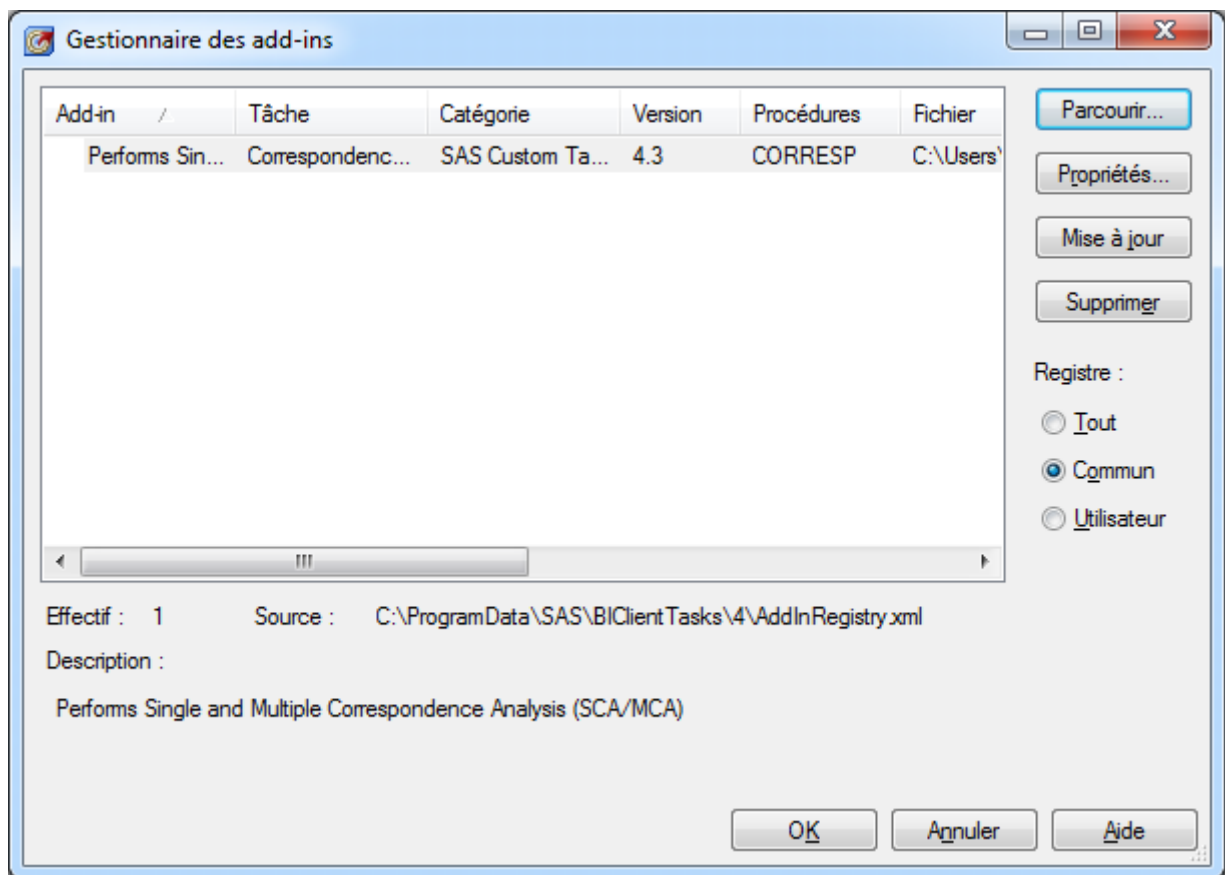
Installation de l'Add-in

Pour réaliser une Analyse en Correspondances Multiples, il vous faudra installer un Add-in (cette add-in, comme tous les add-in, n'est ni traduit, ni supporté, à cause de son statut expérimental), Correspondence Analysis, développé par Audimar Bangi (développeur SAS). La manipulation est simple : il suffit de télécharger le fichier zip en cliquant sur [ce lien](#). Vous y trouverez un fichier texte, deux tables d'exemple et un fichier .dll. La procédure d'installation est la suivante :

1. Copier le fichier SAS.Tasks.CorrespondenceAnalysis.dll dans le répertoire de votre choix, par exemple C:\addinEG
2. Démarrer Enterprise Guide **4.3** (attention : un add-in est spécifique à une version donnée)
3. Cliquer sur « Outils » → « Add-in » → « Gestionnaire des Add-in »



Cliquer sur « Parcourir », sélectionner le fichier .dll et cliquer sur « OK ».



L'add-in est ainsi installé et activé. Il est prêt à être utilisé, à partir du menu « Outils » → « Add-in » → « Correspondence Analysis... » comme nous allons le voir dans les exemples qui vont suivre.

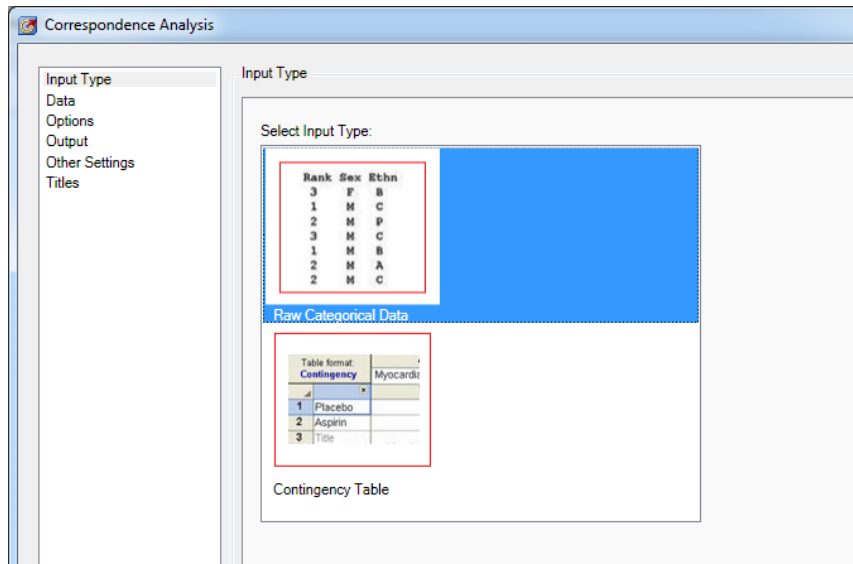
Exemples d'application de l'Add-in d'Analyse en Correspondances Multiples

Nous allons maintenant explorer quelques fonctionnalités offertes par cet outil au travers de plusieurs exemples d'applications, se basant sur les tables d'exemple téléchargées avec l'add-in).

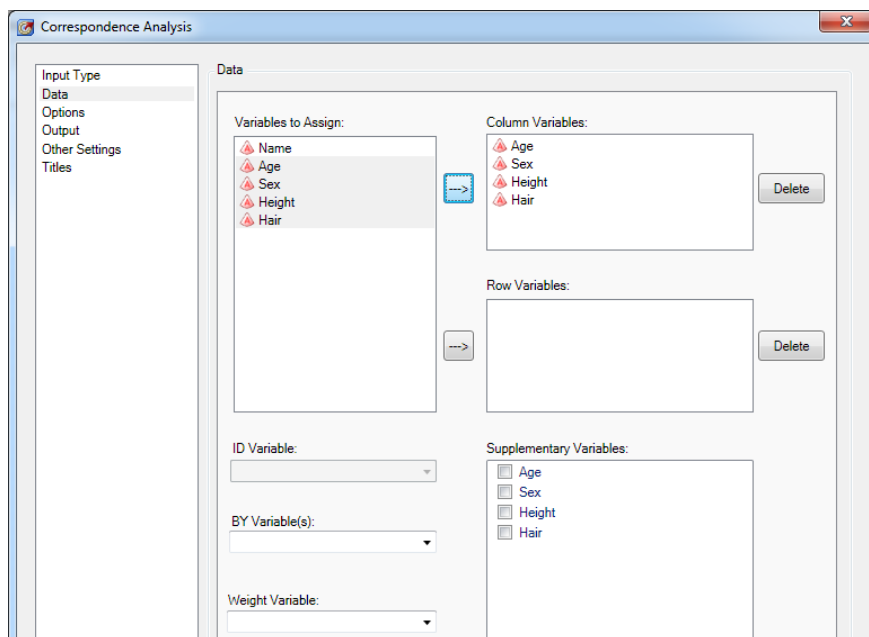
Table avec des variables de classes

La première application se base sur la table categorical data.

- L'ouvrir dans SAS Enterprise Guide 4.3.
- Sélectionner la table et lancer l'analyse grâce au menu (« Outils » → « Add-in » → « Correspondence Analysis... ») et commencer par sélectionner le type d'entrée, par exemple « Raw Categorical Data » (Données catégorielles brutes).

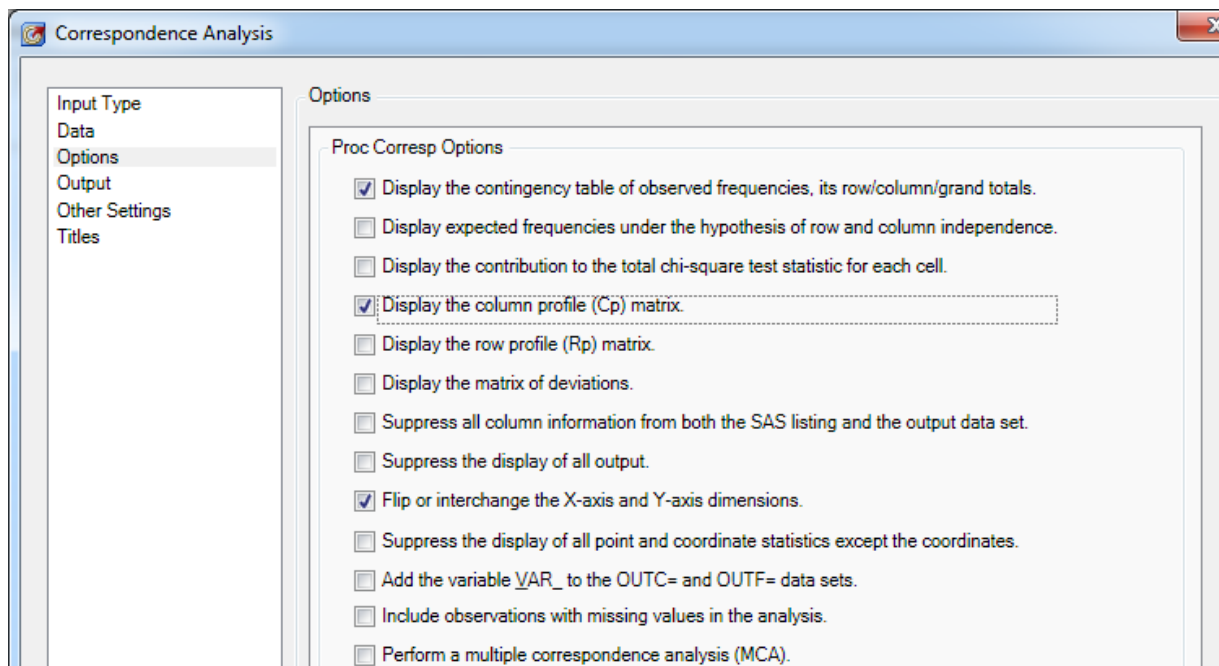


L'exemple s'appuie sur les variables Age, Sex, Height et Hair.

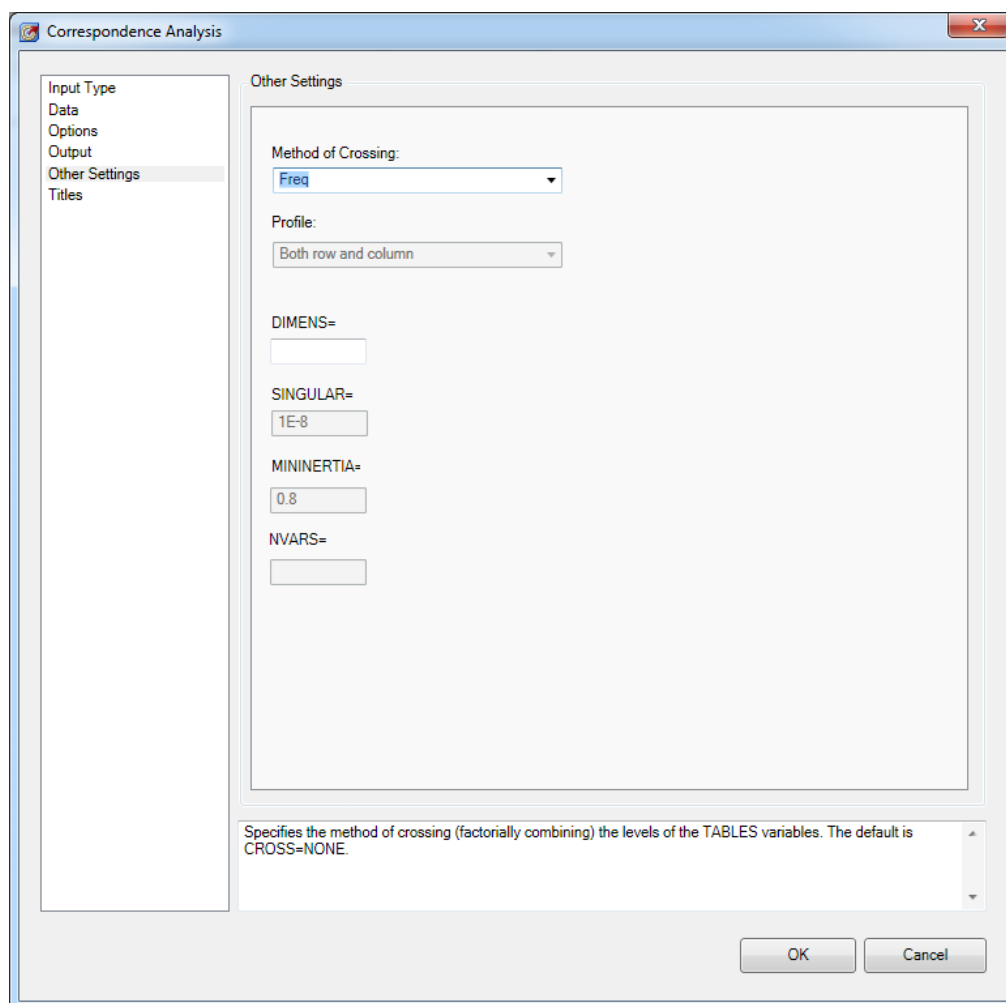


Dans l'onglet "Options", choisir les trois options suivantes (dont voici l'équivalence en Français) :

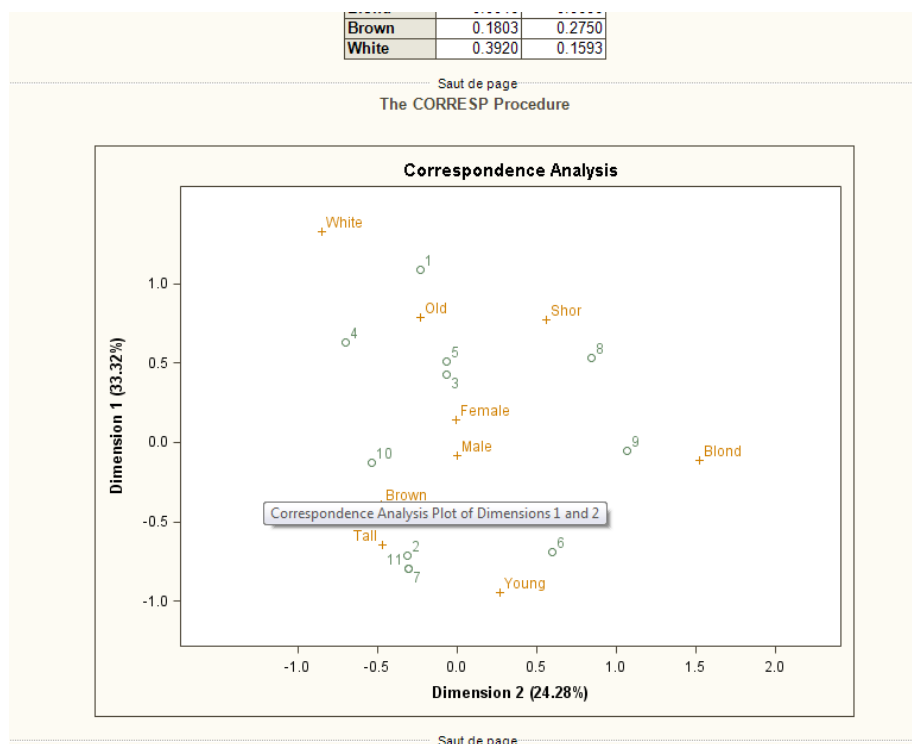
- affichage de la table de contingence des fréquences observées, avec les totaux en lignes, colonnes et global.
- affichage de la matrice de profil de colonne (Cp).
- retourne ou échange les dimensions des axes X et Y.



Dans l'onglet « Other Settings », choisir « Freq » comme méthode de croisement (« Method of Crossing ») et cliquer sur « OK » pour lancer l'analyse.



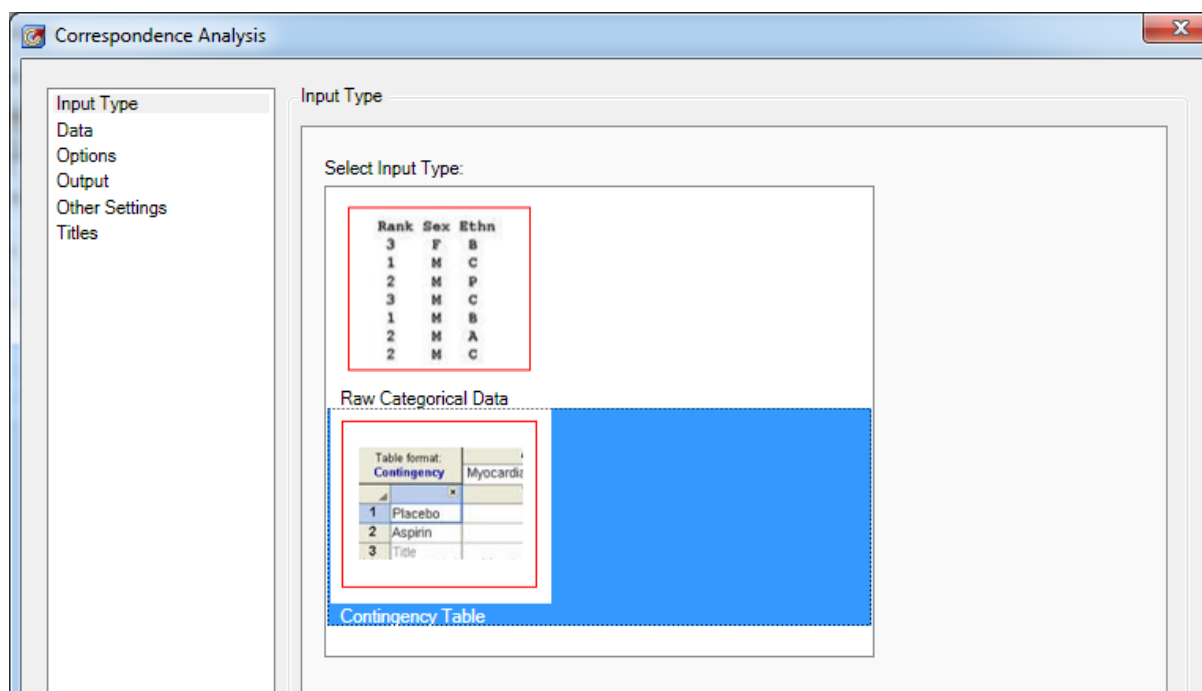
Vous trouverez ci-dessous, un extrait des résultats obtenus en sortie de cette analyse, représentant nos individus et les différentes modalités de nos variables sur un même graphique, où les deux premiers axes sont les deux plus représentatifs.



On note sur le graphique ci-dessus, une sur-représentation des personnes âgées (OLD) avec des cheveux blancs chez les individus classés dans les groupes 1 et 4 alors que le groupe 9 contient une sur-représentation par rapport à la moyenne des blonds.

Application avec une table de contingence

Pour cette seconde analyse, la table contingencytable.sas7bdat sera la table d'étude. Choisir « Table de Contingence » (« Contingency Table ») comme type d'entrée.



Pour cette analyse, sélectionner les six variables du type y19xx comme variables colonnes, « Region » comme « ID variable » (variable d'identification) et w en tant que variable de pondération (« Weight Variable »).

The screenshot shows the 'Correspondence Analysis' window with the 'Data' tab selected. On the left, a sidebar lists 'Input Type', 'Data', 'Options', 'Output', 'Other Settings', and 'Titles'. The main area is divided into several sections:

- Variables to Assign:** A list box containing 'Region', 'y1920', 'y1930', 'y1940', 'y1950', 'y1960', 'y1970', and 'w'. A right-pointing arrow is next to it.
- Column Variables:** A list box containing 'y1920', 'y1930', 'y1940', 'y1950', 'y1960', and 'y1970'. A 'Delete' button is to its right.
- Row Variables:** An empty list box with a 'Delete' button to its right.
- ID Variable:** A dropdown menu currently showing 'Region'.
- BY Variable(s):** An empty dropdown menu.
- Weight Variable:** A dropdown menu currently showing 'w'.
- Supplementary Variables:** A list box containing 'y1920', 'y1930', 'y1940', 'y1950', 'y1960', and 'y1970'.

Le paramétrage des options est identique à celui de la première analyse. Cliquer sur « OK ».

The screenshot shows the 'Options' window for the Correspondence Analysis. The 'Input Type' sidebar on the left has 'Options' selected. The main area is titled 'Proc Corresp Options' and contains a list of checkboxes:

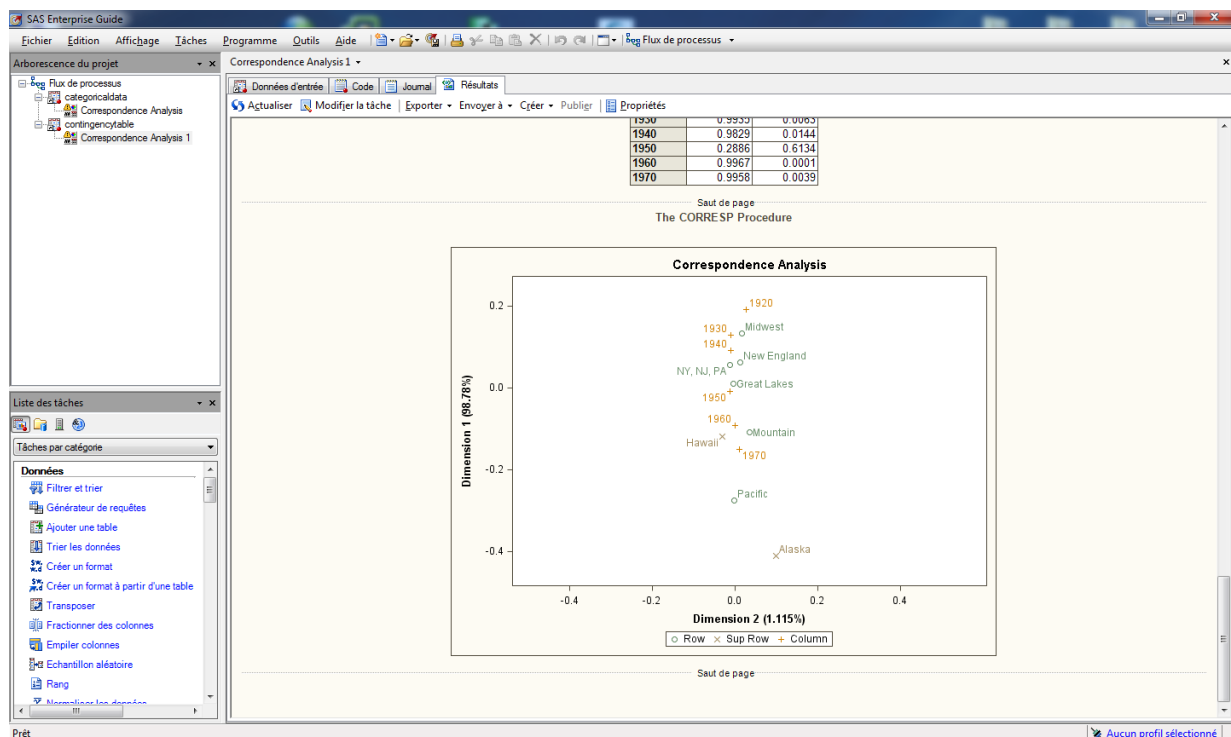
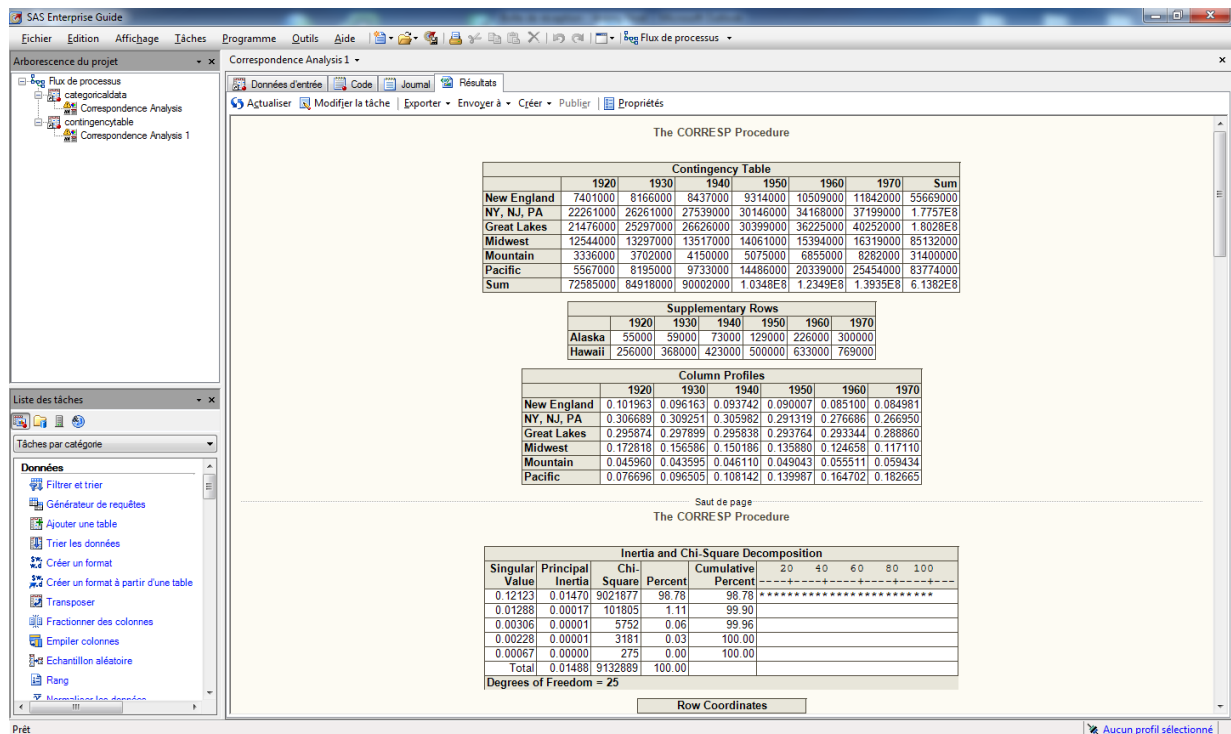
- ☒ Display the contingency table of observed frequencies, its row/column/grand totals.
- ☐ Display expected frequencies under the hypothesis of row and column independence.
- ☐ Display the contribution to the total chi-square test statistic for each cell.
- ☒ Display the column profile (Cp) matrix.
- ☐ Display the row profile (Rp) matrix.
- ☐ Display the matrix of deviations.
- ☐ Suppress all column information from both the SAS listing and the output data set.
- ☐ Suppress the display of all output.
- ☒ Flip or interchange the X-axis and Y-axis dimensions.
- ☐ Suppress the display of all point and coordinate statistics except the coordinates.
- ☐ Add the variable VAR_ to the OUTC= and OUTF= data sets.
- ☐ Include observations with missing values in the analysis.
- ☐ Perform a multiple correspondence analysis (MCA).

Below this is the 'MCA Options' section with three unchecked checkboxes:

- ☐ Display adjusted inertias using Benzecri's method.
- ☐ Display adjusted inertias using Greenacre's method.
- ☐ Display unadjusted inertias.

At the bottom, there is a text box containing 'Flips or interchanges the X-axis and Y-axis dimensions.' and two buttons: 'OK' and 'Cancel'.

Les deux captures d'écrans ci-dessous montrent une partie des résultats (le début et la fin). Ils mettent en évidence que la quasi-totalité de l'inertie du nuage (98.78%) est captée par le premier axe d'analyse, que ce soit à travers le tableau ou le graphique.



Le code SAS traditionnel

Dans ce nouvel Add-in, comme derrière chacune des tâches d'Enterprise Guide, du code SAS est créé et compilé à chaque exécution. La procédure utilisée pour la mise en œuvre de cette analyse est la PROC CORRESP, dont la syntaxe va être détaillée dans cette section. La structure générale est la suivante :


```

PROC CORRESP <options> ;
  TABLES <row-variables,> column-variables ;
  VAR variables ;
  BY variables ;
  ID variable ;
  SUPPLEMENTARY variables ;
  WEIGHT variable ;
run ;

```

Attention les instructions TABLES et VAR sont ressemblantes, il faut en spécifier l'une ou l'autre mais en aucun cas les deux en même temps.

Dans l'instruction PROC CORRESP, en plus de la table d'entrée et des éventuelles tables de sorties (coordonnées et fréquences), il faut, si l'on souhaite faire une ACM, rajouter l'option MCA. Cette option requiert la présence d'un tableau de Burt en entrée. Il en ressort que ces études très francophones sont reprises dans SAS, au travers d'options, puisque ce sont des analyses plus proches de l'exploratoire que du protocole (raison de la présence des ACM dans JMP® ou SAS/IML® Studio, plutôt qu'en standard dans SAS Enterprise Guide).

Les deux exemples d'applications vues précédemment utilisent les mêmes options qui sont :

- observed : affiche une table de contingence des fréquences observées
- cp : affiche la matrice des profils colonnes
- binary: permet la création d'une table binaire facilement
- outc : spécifie le nom de la table de sortie contenant les coordonnées
- outf: spécifie le nom de la table de sortie contenant les fréquences, profils de lignes,...
- print : permet de choisir entre l'affichage en pourcentage ou en fréquence (effectifs, par défaut)
- Plots(flip)=all: produit tous les graphiques appropriés

Vous pouvez vous reporter à [l'aide en ligne](#) sur cette procédure pour plus d'information.

Conclusion

Il est désormais possible de faire une ACM (Analyse en Correspondances Multiples) en mode clic-bouton sous SAS Enterprise Guide, grâce à l'Add-in expérimental que nous venons de présenter. Cependant, il ne faut pas oublier que derrière ce module implémenté, et les fenêtres de SAS Enterprise Guide qui s'en chargent (onglet « Code », « Journal » et « Résultats »), du code SAS tourne en arrière plan avec toutes les spécificités que cela engendre (personnalisation, imbrication,...).

Jérémy NOEL
Consultant Support Clients SAS France